
Dipper Documentation

Tom Conlin, Kent Shefchek, Tim Putman, Matthew Brush, Lilly Wi

Jun 28, 2018

Contents

1 Getting started	3
1.1 Installation	3
1.2 Getting started with Dipper	4
1.3 Notebooks	5
1.4 Downloads	5
1.5 Ingest status	6
1.6 Applications	6
2 Deeper into Dipper	7
2.1 Working with graphs	7
2.2 Working with the models API	8
2.3 Writing ingest with the source API	9
2.4 Testing ingest	11
2.5 Configuring dipper with keys and passwords	12
2.6 Schemas	12
3 For developers	13
3.1 API Docs	13
3.2 Source APIs	63
4 Indices and tables	65
Python Module Index	67

Dipper is a Python package to generate RDF triples from common scientific resources. Dipper includes subpackages and modules to create graphical models of this data, including:

- Models package for generating common sets of triples, including common OWL axioms, complex genotypes, associations, evidence and provenance models.
- Graph package for building graphs with RDFLib or streaming n-triples
- Source package containing fetchers and parsers that interface with remote databases and web services

CHAPTER 1

Getting started

Installing, running, and the basics

1.1 Installation

Dipper requires Python version 3.5 or higher

Install with pip:

```
pip install dipper
```

1.1.1 Development version

The development version can be pulled from GitHub.

```
pip3 install git+git://github.com/monarch-initiative/dipper.git
```

1.1.2 Building locally

```
git clone https://github.com/monarch-initiative/dipper.git
cd dipper
pip install .
```

Alternatively, a subset of source specific requirements may be downloaded. To download the core requirements:

```
pip install -r requirements.txt
```

To download source specific requirements use the requirements/ directory, for example:

```
pip install -r requirements/mgi.txt
```

To download requirements for all sources:

```
pip install -r requirements/all-sources.txt
```

1.2 Getting started with Dipper

This guide assumes you have already installed dipper. If not, then follow the steps in the [Installation](#) section.

1.2.1 Command line

You can run the code by supplying a list of one or more sources on the command line. some examples:

```
dipper-etl.py --sources omim,ncbigene
```

Furthermore, you can check things out by supplying a limit. this will only process the first N number of rows or data elements:

```
dipper-etl.py --sources hpoa --limit 100
```

Other command line parameters are explained if you request help:

```
dipper-etl.py --help
```

1.2.2 Notebooks

We provide [Jupyter Notebooks](#) to illustrate the functionality of the python library. These can also be used interactively.

See the [Notebooks](#) section for more details.

1.2.3 Building models

This code example shows some of the basics of building RDF graphs using the models API:

```
import pandas as pd
from dipper.graph.RDFGraph import RDFGraph
from dipper.models.Model import Model

columns = ['variant', 'variant_label', 'variant_type',
           'phenotype', 'relation', 'source', 'evidence', 'dbxref']

data = [
    ['ClinVarVariant:254143', 'C326F', 'SO:0000694',
     'HP:0000504', 'RO:0002200', 'PMID:12503095', 'ECO:0000220',
     'dbSNP:886037891']
]

# Initialize graph and model
graph = RDFGraph()
```

(continues on next page)

(continued from previous page)

```

model = Model(graph)

# Read file
dataframe = pd.DataFrame(data=data, columns=columns)

for index, row in dataframe.iterrows():
    # Add the triple ClinVarVariant:254143 RO:0002200 HP:0000504
    # RO:0002200 is the has_phenotype relation
    # HP:0000748 is the phenotype 'Inappropriate laughter'
    model.addTriple(row['variant'], row['relation'], row['phenotype'])

    # The addLabel method adds a label using the rdfs:label relation
    model.addLabel(row['variant'], row['variant_label'])

    # addType makes the variant an individual of a class,
    # in this case SO:0000694 'SNP'
    model.addType(row['variant'], row['variant_type'])

    # addXref uses the relation OIO:hasDbXref
    model.addXref(row['variant'], row['dbxref'])

    # Serialize the graph as turtle
    print(graph.serialize(format='turtle').decode("utf-8"))

```

For more information see the [Working with the models API](#) section.

1.3 Notebooks

1.3.1 Jupyter notebook examples

We use Jupyter Notebooks

Available tutorials include:

- Building graphs with the model API
- Querying IMPC evidence and provenance

1.3.2 Running jupyter locally

Follow the instructions for installing from GitHub in [Installation](#). Then start a notebook browser with:

```

pip install jupyter
PYTHONPATH=. jupyter notebook ./docs/notebooks

```

1.4 Downloads

1.4.1 RDF

The dipper output is quality checked and released on a regular basis. The latest release can be found here:

- <https://archive.monarchinitiative.org/latest/ttl/>

The output from our development branch are made available here (may contain errors):

- <https://data.monarchinitiative.org/ttl/>

1.4.2 TSV

TSV downloads for common queries can be found here:

- <https://archive.monarchinitiative.org/latest/tsv/>

1.4.3 Neo4J

A dump of our Neo4J database that includes the output from dipper plus various ontologies:

- <https://archive.monarchinitiative.org/latest/scigraph.tgz>

A public version can be accessed via the SciGraph REST API:

- <https://scigraph-data.monarchinitiative.org/scigraph/docs/>

1.5 Ingest status

We use Jenkins to periodically build each source. A dashboard containing the current status of each ingest can be found here:

- <http://ci.monarchinitiative.org/view/dipper/>

1.6 Applications

1.6.1 BioLink API

A portion of BioLink is driven by the Dipper/SciGraph ETL pipeline:

- <https://api.monarchinitiative.org/api/>

1.6.2 Monarch Initiative

The Monarch application is powered in part by Dipper:

- <https://monarchinitiative.org/>

1.6.3 Owlsim

Annotations loaded into Owlsim are from the Dipper/SciGraph pipeline:

- <http://owlsim3.monarchinitiative.org/api/docs/>

CHAPTER 2

Deeper into Dipper

A look into the structure of the codebase and how to write ingests

2.1 Working with graphs

The Dipper graph package provides two graph implementations, a `RDFGraph` which is an extension of the `RDFLib1` `Graph2`, and a `StreamedGraph` which prints triples to standard out in the ntriples format.

2.1.1 RDFGraphs

The `RDFGraph` class reads the `curie_map.yaml` file and converts strings formatted as curies to `RDFLib` `URIRefs`. Triples are added via the `addTriple` method, for example:

```
from dipper.graph.RDFGraph import RDFGraph

graph = RDFGraph()
graph.addTriple('foaf:John', 'foaf:knows', 'foaf:Joseph')
```

The graph can then be serialized in a variety of formats using `RDFLib3`:

```
from dipper.graph.RDFGraph import RDFGraph

graph = RDFGraph()
graph.addTriple('foaf:John', 'foaf:knows', 'foaf:Joseph')
print(graph.serialize(format='turtle').decode("utf-8"))

# Or write to file
graph.serialize(destination="/path/to/output.ttl", format='turtle')
```

¹ `RDFLib`: <http://rdflib.readthedocs.io/en/stable/>

² `RDFLib Graphs`: <https://rdflib.readthedocs.io/en/stable/apidocs/rdflib.html#graph-module>

³ `RDFLib Serializing`: <http://rdflib.readthedocs.io/en/stable/apidocs/rdflib.html#rdflib.graph.Graph.serialize>

Prints:

```
@prefix OBO: <http://purl.obolibrary.org/obo/> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix xml: <http://www.w3.org/XML/1998/namespace> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .

foaf:John foaf:knows foaf:Joseph .
```

When an object is a literal, set the object_is_literal param to True

```
from dipper.graph.RDFGraph import RDFGraph

graph = RDFGraph()
graph.addTriple('foaf:John', 'rdfs:label', 'John', object_is_literal=True)
```

Literal types can also be passed into the method:

```
from dipper.graph.RDFGraph import RDFGraph

graph = RDFGraph()
graph.addTriple(
    'foaf:John', 'foaf:age', 12,
    object_is_literal=True, literal_type="xsd:integer"
)
```

2.1.2 StreamedGraphs

StreamedGraphs print triples as they are processed by the addTriple method. This is useful for large sources where. The output should be sorted and unqualified as there is no checking for duplicate triples. For example:

```
from dipper.graph.StreamedGraph import StreamedGraph

graph = StreamedGraph()
graph.addTriple('foaf:John', 'foaf:knows', 'foaf:Joseph')
```

Prints:

```
<http://xmlns.com/foaf/0.1/John> <http://xmlns.com/foaf/0.1/knows> <http://xmlns.com/
˓→foaf/0.1/Joseph> .
```

2.1.3 References

2.2 Working with the models API

The models package provides classes for building common sets of triples based on our modeling patterns.

For an example see the notebook on this topic: [Building graphs with the model API](#)

2.2.1 Basics

The models class provides methods for building common RDF and OWL statements

For a list of methods, [see the API docs](#).

2.2.2 Building associations

We use the RDF Reification¹ pattern to create ternary statements, for example, adding frequency data to phenotype to disease associations. We utilize the Open Biomedical Association ontology² to reify statements, and the SEPIO ontology to add evidence and provenance.

For a list of classes and methods, [see the API docs](#).

2.2.3 Building genotypes

We use the GENO ontology⁴ to build complex genotypes and their parts.

For a list of methods, [see the API docs](#).

GENO docs: The Genotype Ontology (GENO)

2.2.4 Building complex evidence and provenance graphs

We use the SEPIO ontology to build complex evidence and provenance. For an example see the IMPC source ingest.

For a list of methods, see the API docs for `evidence` and `provenance`.

SEPIO docs: The Scientific Evidence and Provenance Information Ontology (SEPIO)

2.2.5 References

2.3 Writing ingests with the source API

2.3.1 Overview

Although not required to write an ingest, we have provided a source parent class that can be extended to leverage reusable functionality in each ingest.

To create a new ingest using this method, first extend the Source class.

If the source contains flat files, include a files dictionary with this structure:

```
files = {
    'somekey': {
        'file': 'filename.tsv',
        'url': 'http://example.org/filename.tsv'
    },
    ...
}
```

¹ RDF Reification: <https://www.w3.org/TR/rdf-primer/#reification>

² OBAN: <https://github.com/EBISPORT/OBAN>

⁴ GENO: <https://github.com/monarch-initiative/GENO-ontology>

For example:

```
from dipper.sources.Source import Source

class TPO(Source):
    """
    The ToxicophenomicOmicsDB contains data on ...
    """

    files = {
        'genes': {
            'file': 'genes.tsv',
            'url': 'http://example.org/genes.tsv'
        }
    }
```

2.3.2 Initializing the class

Each source class takes a graph_type (string) and are_bnodes_skolemized (boolean) parameters. These parameters are used to initialize a graph object in the Source constructor.

Note: In the future this may be adjusted so that a graph object is passed into each source.

For example:

```
def __init__(self, graph_type, are_bnodes_skolemized):
    super().__init__(graph_type, are_bnodes_skolemized, 'TPO')
```

2.3.3 Writing the fetcher

This method is intended to fetch data from the remote locations (if it is newer than the local copy).

Extend the parent `fetch` function. If a the remote file has already been downloaded. The fetch method checks the remote headers to see if it has been updated. For sources not served over HTTP, this method may need to be overriden, for example in [Bgee](#).

For example:

```
def fetch(self, is_dl_forced=False):
    """
    Fetches files from TPO

    :param is_dl_forced (bool): Force download
    :return None
    """
    self.get_files(is_dl_forced)
```

2.3.4 Writing the parser

Typically these are written by looping through the series of files that were obtained by the fetch method. The goal is to process each file minimally, adding classes and individuals as necessary, and adding triples to the sources' graph.

For example:

```

def parse(self, limit=None):
    """
    Parses genes from TPO

    :param limit (int, optional) limit the number of rows processed
    :return None
    """

    if limit is not None:
        logger.info("Only parsing first %d rows", limit)

    # Open file
    fh = open('/'.join((self.rawdir, self.files['genes']['file'])), 'r')
    # Parse file
    self._add_gene_toxicology(fh, limit)
    # Close file
    fh.close()

```

Considerations when writing a parser

There are certain conventions that we follow when parsing data:

1. Genes are a special case of genomic feature that are added as (OWL) Classes. But all other genomic features are added as individuals of an owl class.
2. If a source references an external identifier then, assume that it has been processed in another source script, and only add the identifier (but not the label) to it within this source's file. This will help prevent label collisions related to slightly different versions of the source data when integrating downstream.
3. You can instantiate a class or individual as many times as you want; they will get merged in the graph and will only show up once in the resulting output.

2.4 Testing ingests

2.4.1 Unit tests

Unit style tests can be achieved by mocking source classes (or specific functions) and testing single functions. The `test_graph_equality` function can be used to test graph equality by supplying a string formatted as headless (no prefixes) turtle and a graph object. Most dipper methods are not pure functions, and rely on side effects to a graph object. Therefore it is best to clear the graph object before any testing logic, eg:

```

from dipper.utils.TestTools import TestUtils

source.graph = RDFGraph(True) # Reset graph
test_util = TestUtils()
source.run_some_function()
expected_triples = """
    foaf:person1 foaf:knows foaf:person2 .
"""
self.assertTrue(test_util.test_graph_equality(
    expected_triples, source.graph))

```

2.4.2 Integration tests

Integration tests can be executed by generating a file that contains a subset of a source’s data in the same format, and running it through the `source.parse()` method, serializing the graph, and then testing this file in some other piece of code or database.

You may see testing code within source classes, but these tests will be deleted or refactored and moved to the test directory.

2.5 Configuring dipper with keys and passwords

Add private configuration parameters into your private `conf.json` file. Examples of items to put into the config include:

- database connection parameters (in the “`dbauth`” object)
- ftp login credentials
- api keys (in the “`keys`” object)

These are organized such that within any object (`dbauth`, `keys`, etc), they are keyed again by the source’s name.

Here is an example:

```
{  
  "keys": {  
    "omim" : "foo",  
  },  
  "dbauth" : {  
    "mgi" : {  
      "user" : "bar",  
      "password" : "baz"  
    }  
  }  
}
```

This file must be placed in the `dipper` package directory and named `conf.json`. If building locally this is in the `dipper/dipper/` directory. If installed with pip this will be in `path/to/env/lib/python3.x/site-packages/dipper/` directory.

2.6 Schemas

Although RDF is inherently schemaless, we aim to construct consistent models across sources. This allows us to build source agnostic queries and bridge data across sources.

The `dipper` schemas are documented as directed graphs. Examples can be found in the [ingest artifacts repo](#).

Some ontologies contain documentation on how to describe data using the classes and properties defined in the ontology:

- The Scientific Evidence and Provenance Information Ontology (SEPIO)
- The Genotype Ontology (GENO)

While not yet implemented, in the future we plan on defining our schemas and constraints using the [BioLink model specification](#).

The cypher queries that we use to cache inferred and direct relationships between entities are [stored in GitHub](#).

CHAPTER 3

For developers

3.1 API Docs

3.1.1 dipper package

Subpackages

dipper.graph package

Submodules

dipper.graph.Graph module

```
class dipper.graph.Graph
Bases: object

addTriple (subject_id, predicate_id, object_id, object_is_literal, literal_type)
serialize (subject_iri, predicate_iri, obj, object_is_literal, literal_type)
skolemizeBlankNode (curie)
```

dipper.graph.RDFGraph module

```
class dipper.graph.RDFGraph (are_bnodes_skized=True, identifier=None)
Bases: rdflib.graph.ConjunctiveGraph, dipper.graph.Graph
```

Extends RDFLibs ConjunctiveGraph The goal of this class is wrap the creation of triples and manage creation of URIRef, Bnodes, and literals from an input curie

```
addTriple (subject_id, predicate_id, obj, object_is_literal=False, literal_type=None)
bind_all_namespaces ()
```

```
curie_util = <dipper.utils.CurieUtil.CurieUtil object>
skolemizeBlankNode (curie)
```

dipper.graph.StreamedGraph module

```
class dipper.graph.StreamedGraph.StreamedGraph (are_bnodes_skized=True,
                                                file_handle=None, fmt='nt')
Bases: dipper.graph.Graph.Graph

Stream rdf triples to file or stdout Assumes a downstream process will sort then uniquify triples
Theoretically could support both ntriple, rdfxml formats, for now just support nt
addTriple (subject_id, predicate_id, object_id, object_is_literal=False, literal_type=None)
curie_util = <dipper.utils.CurieUtil.CurieUtil object>
serialize (subject_iri, predicate_iri, obj, object_is_literal=False, literal_type=None)
skolemizeBlankNode (curie)
```

dipper.models package

Subpackages

dipper.models.assoc package

Submodules

dipper.models.assoc.Association module

```
class dipper.models.assoc.Association.Assoc (graph, definedby, sub=None, obj=None,
                                              pred=None)
Bases: object
```

A base class for OBAN (Monarch)-style associations, to enable attribution of source and evidence on statements.

add_association_to_graph()

add_date (date)

add_evidence (identifier)

Add an evidence code to the association object (maintained as a list) :param identifier:

Returns

add_predicate_object (predicate, object_node, object_type=None, datatype=None)

add_provenance (identifier)

add_source (identifier)

Add a source identifier (such as publication id) to the association object (maintained as a list) TODO we need to greatly expand this function!

Parameters identifier –

Returns

```

annotation_properties = {'consider': 'OIO:consider', 'definition': 'IAO:0000115', 'h
assoc_types = {'association': 'OBAN:association'}
datatype_properties = {'created_on': 'pav:createdOn', 'has_measurement': 'IAO:000000
get_association_id()
get_properties()

static make_association_id(definedby, subject, predicate, object, attributes=None)

```

A method to create unique identifiers for OBAN-style associations, based on all the parts of the association. If any of the items is empty or None, it will convert it to blank. It effectively digests the string of concatenated values. Subclasses of Assoc can submit an additional array of attributes that will be appended to the ID.

Parameters

- **definedby** – The (data) resource that provided the annotation
- **subject** –
- **predicate** –
- **object** –
- **attributes** –

Returns

```

object_properties = {'causes_or_contributes': 'RO:0003302', 'expressed_in': 'RO:0002
properties = {'causes_or_contributes': 'RO:0003302', 'consider': 'OIO:consider', 'cr
set_association_id(assoc_id=None)

```

This will set the association ID based on the internal parts of the association. To be used in cases where an external association identifier should be used.

Parameters **assoc_id** –

Returns

```

set_description(description)
set_object(identifier)
set_relationship(identifier)
set_score(score, unit=None, score_type=None)
set_subject(identifier)

```

dipper.models.assoc.Chem2DiseaseAssoc module

```

class dipper.models.assoc.Chem2DiseaseAssoc.Chem2DiseaseAssoc(graph,      de-
                                                               finedby, chem_id,
                                                               phenotype_id,
                                                               rel_id=None)

```

Bases: *dipper.models.assoc.Association.Assoc*

Attributes: assoc_id (str): Association Curie (Prefix:ID) chem_id (str): Chemical Curie phenotype_id (str): Phenotype Curie pub_list (str,list): One or more publication curies rel (str): Property relating assoc_id and chem_id evidence (str): Evidence curie

```
make_c2p_assoc_id()
```

`set_association_id(assoc_id=None)`

This will set the association ID based on the internal parts of the association. To be used in cases where an external association identifier should be used.

Parameters `assoc_id` –

Returns

dipper.models.assoc.D2PAssoc module

`class dipper.models.assoc.D2PAssoc(graph, definedby, disease_id, phenotype_id, onset=None, frequency=None, rel=None)`

Bases: `dipper.models.assoc.Association.Assoc`

A specific association class for defining Disease-to-Phenotype relationships. This assumes that a graph is created outside of this class, and nodes get added. By default, an association will assume the “has_phenotype” relationship, unless otherwise specified.

`add_association_to_graph()`

The reified relationship between a disease and a phenotype is decorated with some provenance information. This makes the assumption that both the disease and phenotype are classes.

Parameters `g` –

Returns

`d2p_object_properties = {'frequency': ':frequencyOfPhenotype', 'onset': ':onset'}`

`make_d2p_id()`

Make an association id for phenotypic associations with disease that is defined by: source of association + disease + relationship + phenotype + onset + frequency

Returns

`set_association_id(assoc_id=None)`

This will set the association ID based on the internal parts of the association. To be used in cases where an external association identifier should be used.

Parameters `assoc_id` –

Returns

dipper.models.assoc.DispositionAssoc module

`class dipper.models.assocDispositionAssoc(graph, definedby, entity_id, heritability_id)`

Bases: `dipper.models.assoc.Association.Assoc`

A specific Association model for Heritability annotations. These are to be used between diseases and a heritability disposition.

dipper.models.assoc.G2PAssoc module

`class dipper.models.assoc.G2PAssoc(graph, definedby, entity_id, phenotype_id, rel=None)`

Bases: `dipper.models.assoc.Association.Assoc`

A specific association class for defining Genotype-to-Phenotype relationships. This assumes that a graph is created outside of this class, and nodes get added. By default, an association will assume the “has_phenotype” relationship, unless otherwise specified. Note that genotypes are expected to be created and defined outside of this association, most likely by calling methods in the Genotype() class.

add_association_to_graph()

Overrides Association by including bnode support

The reified relationship between a genotype (or any genotype part) and a phenotype is decorated with some provenance information. This makes the assumption that both the genotype and phenotype are classes.

currently hardcoded to map the annotation to the monarch namespace :param g: :return:

```
g2p_types = {'developmental_process': 'GO:0032502'}
```

make_g2p_id()

Make an association id for phenotypic associations that is defined by: source of association + (Annotation subject) + relationship + phenotype/disease + environment + start stage + end stage

Returns

set_association_id(assoc_id=None)

This will set the association ID based on the internal parts of the association. To be used in cases where an external association identifier should be used.

Parameters assoc_id –

Returns

```
set_environment(environment_id)
```

```
set_stage(start_stage_id, end_stage_id)
```

dipper.models.assoc.InteractionAssoc module

```
class dipper.models.assoc.InteractionAssoc.InteractionAssoc(graph, definedby,  
subj, obj, rel=None)
```

Bases: *dipper.models.assoc.Association.Assoc*

```
interaction_object_properties = {'colocalizes_with': 'RO:0002325', 'genetically_interacts_with': 'RO:0002326'}
```

dipper.models.assoc.OrthologyAssoc module

```
class dipper.models.assoc.OrthologyAssoc.OrthologyAssoc(graph, definedby, gene1,  
gene2, rel=None)
```

Bases: *dipper.models.assoc.Association.Assoc*

add_gene_family_to_graph(family_id)

Make an association between a group of genes and some grouping class. We make the assumption that the genes in the association are part of the supplied family_id, and that the genes have already been declared as classes elsewhere. The family_id is added as an individual of type DATA:gene_family.

Triples: <family_id> a DATA:gene_family <family_id> RO:has_member <gene1> <family_id>
RO:has_member <gene2>

Parameters

- **family_id** –
- **g** – the graph to modify

Returns

```
ortho_rel = {'has_member': 'RO:0002351', 'homologous': 'RO:HOM0000019', 'in_paralogous': 'RO:0002352'}  
terms = {'gene_family': 'DATA:3148'}
```

Submodules

dipper.models.Dataset module

```
class dipper.models.Dataset(Dataset(identifier, title, url, description=None, license_url=None, data_rights=None, graph_type=None, file_handle=None))
```

Bases: object

this will produce the metadata about a dataset following the example laid out here: http://htmlpreview.github.io/?https://github.com/joejimbo/HCLSDatasetDescriptions/blob/master/Overview.html#appendix_1 (mind the wrap)

`getGraph()`

`get_license()`

`setFileAccessUrl(url, is_object_literal=False)`

`setVersion(date_issued, version_id=None)`

Legacy function... should use the other set_* for version and date

as of 2016-10-20 used in:

dipper/sources/HPOAnnotations.py 139: dipper/sources/CTD.py 99: dipper/sources/BioGrid.py 100:
dipper/sources/MGI.py 255: dipper/sources/EOM.py 93: dipper/sources/Coriell.py 200: dipper/sources/MMRRC.py 77:

TODO set as deprecated

Parameters

- `date_issued` –
- `version_id` –

Returns

`set_citation(citation_id)`

`set_date_issued(date_issued)`

`set_license(license)`

`set_version_by_date(date_issued=None)`

This will set the version by the date supplied, the date already stored in the dataset description, or by the download date (today) :param date_issued: :return:

`set_version_by_num(version_num)`

dipper.models.Environment module

```
class dipper.models.Environment(Environment(graph))
```

Bases: object

These methods provide convenient methods to add items related to an experimental environment and it's parts to a supplied graph.

This is a stub ready for expansion.

```
addComponentAttributes (component_id, entity_id, value=None, unit=None)
addComponentToEnvironment (env_id, component_id)
addEnvironment (env_id, env_label, env_type=None, env_description=None)
addEnvironmentalCondition (cond_id, cond_label, cond_type=None, cond_description=None)
annotation_properties = {}
environment_parts = {'crispr_reagent': 'REO:crispr_TBD', 'environmental_condition':
object_properties = {'has_part': 'BFO:0000051'}
properties = {'has_part': 'BFO:0000051'}
```

dipper.models.Evidence module

```
class dipper.models.Evidence.Evidence (graph, association)
Bases: object
```

To model evidence as the basis for an association. This encompasses:

- **measurements taken from the lab, and their significance.** these can be derived from papers or other agents.
- papers

>1 measurement may result from an assay, each of which may have it's own significance

```
add_data_individual (data_curie, label=None, ind_type=None)
Add data individual :param data_curie: str either curie formatted or long string,
long strings will be converted to bnodes
```

Parameters

- **type** – str curie
- **label** – str

Returns

None

```
add_evidence (evidence_line, ev_type=None, label=None)
Add line of evidence node to association id
```

Parameters

- **assoc_id** – curie or iri, association id
- **evidence_line** – curie or iri, evidence line

Returns

```
add_source (evidence_line, source, label=None, src_type=None)
Applies the triples: <evidence> <dc:source> <source> <source> <rdf:type> <type> <source> <rdfs:label>
“label”
```

TODO this should belong in a higher level class :param evidence_line: str curie :param source: str source as curie :param label: optional, str type as curie :param type: optional, str type as curie :return: None

add_supporting_data (*evidence_line, measurement_dict*)

Add supporting data :param evidence_line: :param data_object: dict, where keys are curies or iris and values are measurement values for example:

```
{ “_:1234” : “1.53E07” “_:4567”: “20.25”
}
```

Note: assumes measurements are RDF:type ‘ed elsewhere :return: None

add_supporting_evidence (*evidence_line, type=None, label=None*)

Add supporting line of evidence node to association id

Parameters

- **assoc_id** – curie or iri, association id
- **evidence_line** – curie or iri, evidence line

Returns

None

add_supporting_publication (*evidence_line, publication, label=None, pub_type=None*)

<evidence> <SEPIO:0000124> <source> <source> <rdf:type> <type> <source> <rdfs:label> “label” :param evidence_line: str curie :param publication: str curie :param label: optional, str type as curie :param type: optional, str type as curie :return:

```
data_property = { 'has_measurement': 'IAO:0000004', 'has_value': 'STATO:0000129'}
data_types = { 'count': 'SIO:000794', 'odds_ratio': 'STATO:0000182', 'proportional_re...
evidence_types = { 'assay': 'OBI:0000070', 'blood test evidence': 'ECO:0001016', 'eff...
object_properties = { 'has_evidence': 'SEPIO:0000006', 'has_significance': 'STATO:has_...
```

dipper.models.Family module

class dipper.models.Family(*graph*)

Bases: object

Model mereological/part whole relationships

Although these relations are more abstract, we often use them to model family relationships (proteins, humans, etc.) The naming of this class may change in the future to better reflect the meaning of the relations it is modeling

addMember (*group_id, member_id*)

addMemberOf (*member_id, group_id*)

```
object_properties = { 'has_member': 'RO:0002351', 'member_of': 'RO:0002350' }
```

dipper.models.GenomicFeature module

class dipper.models.GenomicFeature.FEATURE(*graph, feature_id=None, label=None, feature_type=None, description=None*)

Bases: object

Dealing with genomic features here. By default they are all faldo:Regions. We use SO for typing genomic features. At the moment, RO:has_subsequence is the default relationship between the regions, but this should be tested/verified.

TODO: the graph additions are in the addXToFeature functions, but should be separated. TODO: this will need to be extended to properly deal with fuzzy positions in faldo.

addFeatureEndLocation (*coordinate, reference_id, strand=None, position_types=None*)

Adds the coordinate details for the end of this feature :param coordinate: :param reference_id: :param strand:

Returns

addFeatureProperty (*property_type, property*)

addFeatureStartLocation (*coordinate, reference_id, strand=None, position_types=None*)

Adds coordinate details for the start of this feature. :param coordinate: :param reference_id: :param strand: :param position_types:

Returns

addFeatureToGraph (*add_region=True, region_id=None, feature_as_class=False*)

We make the assumption here that all features are instances. The features are located on a region, which begins and ends with faldo:Position. The feature locations leverage the Faldo model, which has a general structure like: Triples: feature_id a feature_type (individual) faldo:location region_id region_id a faldo:region faldo:begin start_position faldo:end end_position start_position a (any of: faldo:(((Both|Plus|Minus)Strand)|Exact)Position) faldo:position Integer(numeric position) faldo:reference reference_id end_position a (any of: faldo:(((Both|Plus|Minus)Strand)|Exact)Position) faldo:position Integer(numeric position) faldo:reference reference_id

Parameters graph –

Returns

addPositionToGraph (*reference_id, position, position_types=None, strand=None*)

Add the positional information to the graph, following the faldo model. We assume that if the strand is None, we give it a generic “Position” only. Triples: my_position a (any of: faldo:(((Both|Plus|Minus)Strand)|Exact)Position) faldo:position Integer(numeric position) faldo:reference reference_id

Parameters

- **graph –**
- **reference_id –**
- **position –**
- **position_types –**
- **strand –**

Returns Identifier of the position created

addRegionPositionToGraph (*region_id, begin_position_id, end_position_id*)

addSubsequenceOfFeature (*parentid*)

This will add reciprocal triples like: feature is_subsequence_of parent parent has_subsequence feature :param graph: :param parentid:

Returns

`addTaxonToFeature (taxonid)`

Given the taxon id, this will add the following triple: feature in_taxon taxonid :param graph: :param taxonid: :return:

```
annotation_properties = {}
```

```
data_properties = {'position': 'faldo:position'}
```

```
object_properties = {'begin': 'faldo:begin', 'downstream_of_sequence_of': 'RO:000252'}
```

```
properties = {'begin': 'faldo:begin', 'downstream_of_sequence_of': 'RO:0002529', 'en'}
```

```
types = {'FuzzyPosition': 'faldo:FuzzyPosition', 'Position': 'faldo:Position', 'SNP'}
```

```
dipper.models.GenomicFeature.makeChromID (chrom, reference=None, prefix=None)
```

This will take a chromosome number and a NCBI taxon number, and create a unique identifier for the chromosome. These identifiers are made in the @base space like: Homo sapiens (9606) chr1 ==> :9606chr1 Mus musculus (10090) chrX ==> :10090chrX

Parameters

- **chrom** – the chromosome (preferably without any chr prefix)
- **reference** – the numeric portion of the taxon id

Returns

```
dipper.models.GenomicFeature.makeChromLabel (chrom, reference=None)
```

dipper.models.Genotype module

`class dipper.models.Genotype.Genotype (graph)`

Bases: object

These methods provide convenient methods to add items related to a genotype and its parts to a supplied graph. They follow the patterns set out in GENO <https://github.com/monarch-initiative/GENO-ontology>. For specific sequence features, we use the GenomicFeature class to create them.

`addAffectedLocus (allele_id, gene_id, rel_id=None)`

We make the assumption here that if the relationship is not provided, it is a GENO:is_sequence_variant_instance_of.

Here, the allele should be a variant_locus, not a sequence alteration. :param allele_id: :param gene_id: :param rel_id: :return:

`addAllele (allele_id, allele_label, allele_type=None, allele_description=None)`

Make an allele object. If no allele_type is added, it will default to a geno:allele :param allele_id: curie for allele (required) :param allele_label: label for allele (required) :param allele_type: id for an allele type (optional, recommended SO or GENO class) :param allele_description: a free-text description of the allele :return:

`addAlleleOfGene (allele_id, gene_id, rel_id=None)`

We make the assumption here that if the relationship is not provided, it is a GENO:is_sequence_variant_instance_of.

Here, the allele should be a variant_locus, not a sequence alteration. :param allele_id: :param gene_id: :param rel_id: :return:

`addChromosome (chr, tax_id, tax_label=None, build_id=None, build_label=None)`

If it's just the chromosome, add it as an instance of a SO:chromosome, and add it to the genome. If a build is included, punn the chromosome as a subclass of SO:chromosome, and make the build-specific chromosome an instance of the supplied chr. The chr then becomes part of the build or genome.

addChromosomeClass (*chrom_num*, *taxon_id*, *taxon_label*)

addChromosomeInstance (*chr_num*, *reference_id*, *reference_label*, *chr_type=None*)

Add the supplied chromosome as an instance within the given reference :param *chr_num*: :param *reference_id*: for example, a build id like UCSC:hg19 :param *reference_label*: :param *chr_type*: this is the class that this is an instance of. typically a genome-specific chr

Returns

addConstruct (*construct_id*, *construct_label*, *construct_type=None*, *construct_description=None*)

addDerivesFrom (*child_id*, *parent_id*)

We add a derives_from relationship between the child and parent id. Examples of uses include between: an allele and a construct or strain here, a cell line and it's parent genotype. Adding the parent and child to the graph should happen outside of this function call to ensure graph integrity. :param *child_id*: :param *parent_id*: :return:

addGene (*gene_id*, *gene_label*, *gene_type=None*, *gene_description=None*)

addGeneProduct (*sequence_id*, *product_id*, *product_label=None*, *product_type=None*)

Add gene/variant/allele has_gene_product relationship Can be used to either describe a gene to transcript relationship or gene to protein :param *sequence_id*: :param *product_id*: :param *product_label*: :param *product_type*: :return:

addGeneTargetingReagent (*reagent_id*, *reagent_label*, *reagent_type*, *gene_id*, *description=None*)

Here, a gene-targeting reagent is added. The actual targets of this reagent should be added separately. :param *reagent_id*: :param *reagent_label*: :param *reagent_type*:

Returns

addGeneTargetingReagentToGenotype (*reagent_id*, *genotype_id*)

addGenome (*taxon_id*, *taxon_label=None*)

addGenomicBackground (*background_id*, *background_label*, *background_type=None*, *background_description=None*)

addGenomicBackgroundToGenotype (*background_id*, *genotype_id*, *background_type=None*)

addGenotype (*genotype_id*, *genotype_label*, *genotype_type=None*, *genotype_description=None*)

If a genotype_type is not supplied, we will default to ‘intrinsic_genotype’ :param *genotype_id*: :param *genotype_label*: :param *genotype_type*: :param *genotype_description*: :return:

addMemberOfPopulation (*member_id*, *population_id*)

addParts (*part_id*, *parent_id*, *part_relationship=None*)

This will add a has_part (or subproperty) relationship between a parent_id and the supplied part. By default the relationship will be BFO:has_part, but any relationship could be given here. :param *part_id*: :param *parent_id*: :param *part_relationship*: :return:

addPartsToVSLC (*vslc_id*, *allele1_id*, *allele2_id*, *zygosity_id=None*, *allele1_rel=None*, *allele2_rel=None*)

Here we add the parts to the VSLC. While traditionally alleles (reference or variant loci) are traditionally added, you can add any node (such as sequence_alterations for unlocated variations) to a vslc if they are known to be paired. However, if a sequence_alteration’s loci is unknown, it probably should be added directly to the GVC. :param *vslc_id*: :param *allele1_id*: :param *allele2_id*: :param *zygosity_id*: :param *allele1_rel*: :param *allele2_rel*: :return:

addPolypeptide (*polypeptide_id*, *polypeptide_label=None*, *transcript_id=None*, *polypeptide_type=None*)

Parameters

- **polypeptide_id** –

- **polypeptide_label** -
- **polypeptide_type** -
- **transcript_id** -

Returns

```
addReagentTargetedGene (reagent_id, gene_id, targeted_gene_id=None, tar-  
geted_gene_label=None, description=None)
```

This will create the instance of a gene that is targeted by a molecular reagent (such as a morpholino or rnai). If an instance id is not supplied, we will create it as an anonymous individual which is of the type GENO:reagent_targeted_gene. We will also add the targets relationship between the reagent and gene class.

```
<targeted_gene_id> a GENO:reagent_targeted_gene rdf:label targeted_gene_label dc:description description  
<reagent_id> GENO:targets_instance_of <gene_id>
```

Parameters

- **reagent_id** -
- **gene_id** -
- **targeted_gene_id** -

Returns

```
addReferenceGenome (build_id, build_label, taxon_id)
```

```
addSequenceAlteration (sa_id, sa_label, sa_type=None, sa_description=None)
```

```
addSequenceAlterationToVariantLocus (sa_id, vl_id)
```

```
addSequenceDerivesFrom (child_id, parent_id)
```

```
addTargetedGeneComplement (tgc_id, tgc_label, tgc_type=None, tgc_description=None)
```

```
addTargetedGeneSubregion (tgs_id, tgs_label, tgs_type=None, tgs_description=None)
```

```
addTaxon (taxon_id, genopart_id)
```

The supplied geno part will have the specified taxon added with RO:in_taxon relation. Generally the taxon is associated with a genomic_background, but could be added to any genotype part (including a gene, regulatory element, or sequence alteration). :param taxon_id: :param genopart_id:

Returns

```
addVSLCtoParent (vslc_id, parent_id)
```

The VSLC can either be added to a genotype or to a GVC. The vslc is added as a part of the parent. :param vslc_id: :param parent_id: :return:

```
annotation_properties = {'altered_nucleotide': 'GENO:altered_nucleotide', 'reference_...
```

```
genoparts = {'QTL': 'SO:0000771', 'RNAi_reagent': 'SO:0000337', 'allele': 'GENO:0000...
```

```
makeGenomeID (taxon_id)
```

```
make_experimental_model_with_genotype (genotype_id, genotype_label, taxon_id,  
taxon_label)
```

```
make_variant_locus_label (gene_label, allele_label)
```

```
make_vslc_label (gene_label, allele1_label, allele2_label)
```

Make a Variant Single Locus Complement (VSLC) in monarch-style. :param gene_label: :param allele1_label: :param allele2_label: :return:

```
object_properties = {'derives_from': 'RO:0001000', 'derives_sequence_from_gene': 'GE...
```

```
properties = {'altered_nucleotide': 'GENO:altered_nucleotide', 'derives_from': 'RO:0000001'}
zygosity = {'complex_heterozygous': 'GENO:0000402', 'hemizygous': 'GENO:0000606', 'homozygous': 'GENO:0000607'}
```

dipper.models.Model module

class dipper.models.Model(*graph*)

Bases: object

Utility class to add common triples to a graph (subClassOf, type, label, sameAs)

addBlankNodeAnnotation(*node_id*)

Add an annotation property to the given `node_id` to be a pseudo blank node. This is a monarchism.
:param node_id: :return:

addClassToGraph(*class_id*, *label=None*, *class_type=None*, *description=None*)

Any node added to the graph will get at least 3 triples: *(node, type, owl:Class) and *(node, label, literal(label)) *if a type is added,

then the node will be an OWL:subclassOf that the type

*if a description is provided, it will also get added as a dc:description

Parameters

- **class_id** –
- **label** –
- **class_type** –
- **description** –

Returns

addComment(*subject_id*, *comment*)

addDefinition(*class_id*, *definition*)

addDepiction(*subject_id*, *image_url*)

addDeprecatedClass(*old_id*, *new_ids=None*)

Will mark the oldid as a deprecated class. if one newid is supplied, it will mark it as replaced by. if >1 newid is supplied, it will mark it with consider properties :param old_id: str - the class id to deprecate
:param new_ids: list - the class list that is

the replacement(s) of the old class. Not required.

Returns

addDeprecatedIndividual(*old_id*, *new_ids=None*)

Will mark the oldid as a deprecated individual. if one newid is supplied, it will mark it as replaced by. if >1 newid is supplied, it will mark it with consider properties :param g: :param oldid: the individual id to deprecate :param newids: the individual idlist that is the replacement(s) of

the old individual. Not required.

Returns

addDescription(*subject_id*, *description*)

```
addEquivalentClass (sub, obj)
addIndividualToGraph (ind_id, label, ind_type=None, description=None)
addLabel (subject_id, label)
addOWLPropertyClassRestriction (class_id, property_id, property_value)
addOWLVersionIRI (ontology_id, version_iri)
addOWLVersionInfo (ontology_id, version_info)
addOntologyDeclaration (ontology_id)
addPerson (person_id, person_label=None)
addSameIndividual (sub, obj)
addSubClass (child_id, parent_id)
addSynonym (class_id, synonym, synonym_type=None)
Add the synonym as a property of the class cid. Assume it is an exact synonym, unless otherwise specified
:param g: :param cid: class id :param synonym: the literal synonym label :param synonym_type: the
CURIE of the synonym type (not the URI) :return:
addTriple (subject_id, predicate_id, obj, object_is_literal=False, literal_type=None)
addType (subject_id, subject_type)
addXref (class_id, xref_id, xref_as_literal=False)
annotation_properties = {'clique_leader': 'MONARCH:cliqueLeader', 'comment': 'dc:com-
datatype_properties = {'has_measurement': 'IAO:0000004', 'position': 'faldo:position'
makeLeader (node_id)
Add an annotation property to the given `node_id` to be the clique_leader. This is a monarchism.
:param node_id: :return:
object_properties = {'causally_influences': 'RO:0002566', 'causally_upstream_of_or_wi-
types = {'annotation_property': 'owl:AnnotationProperty', 'class': 'owl:Class', 'dat-
```

dipper.models.Pathway module

```
class dipper.models.Pathway(graph)
Bases: object
```

This provides convenience methods to deal with gene and protein collections in the context of pathways.

```
addComponentToPathway (component_id, pathway_id)
```

This can be used directly when the component is directly involved in the pathway. If a transforming event is performed on the component first, then the addGeneToPathway should be used instead.

Parameters

- `pathway_id` –
- `component_id` –

Returns

```
addGeneToPathway (gene_id, pathway_id)
```

When adding a gene to a pathway, we create an intermediate ‘gene product’ that is involved in the pathway, through a blank node.

gene_id RO:has_gene_product _gene_product _gene_product RO:involved_in pathway_id

Parameters

- **pathway_id** –
- **gene_id** –

Returns

addPathway (*pathway_id*, *pathway_label*, *pathway_type=None*, *pathway_description=None*)

Adds a pathway as a class. If no specific type is specified, it will default to a subclass of “GO:cellular_process” and “PW:pathway”. :param pathway_id: :param pathway_label: :param pathway_type: :param pathway_description: :return:

```
object_properties = {'gene_product_of': 'RO:0002204', 'has_gene_product': 'RO:0002205'}
pathway_parts = {'cellular_process': 'GO:0009987', 'gene_product': 'CHEBI:33695', 'p
properties = {'gene_product_of': 'RO:0002204', 'has_gene_product': 'RO:0002205', 'in
```

dipper.models.Provenance module

class dipper.models.Provenance (*graph*)
Bases: object

To model provenance as the basis for an association. This encompasses:

- Process history leading to a claim being made, including processes through which evidence is evaluated
- Processes through which information used as evidence is created.

Provenance metadata includes accounts of who conducted these processes, what entities participated in them, and when/where they occurred.

```
add_agent_to_graph (agent_id, agent_label, agent_type=None, agent_description=None)
add_assay_to_graph (assay_id, assay_label, assay_type=None, assay_description=None)
add_assertion (assertion, agent, agent_label, date=None)
    Add assertion to graph :param assertion: :param agent: :param evidence_line: :param date: :return: None
add_date_created (prov_type, date)
add_study_measure (study, measure)
add_study_parts (study, study_parts)
add_study_to_measurements (study, measurements)
object_properties = {'asserted_by': 'SEPIO:0000130', 'created_at_location': 'SEPIO:0
provenance_types = {'assay': 'OBI:0000070', 'assertion': 'SEPIO:0000001', 'assertion
```

dipper.models.Reference module

class dipper.models.Reference (*graph*, *ref_id=None*, *ref_type=None*)
Bases: object

To model references for associations (such as journal articles, books, etc.).

By default, references will be typed as “documents”, unless if the type is set otherwise.

If a `short_citation` is set, this will be used for the individual's label. We may wish to subclass this later.

```
addAuthor (author)
addPage (subject_id, page_url)
addRefToGraph ()
addTitle (subject_id, title)
annotation_properties = {'page': 'foaf:page', 'title': 'dc:title'}
ref_types = {'document': 'IAO:0000310', 'journal_article': 'IAO:0000013', 'person':
setAuthorList (author_list)

    Parameters author_list – Array of authors

    Returns

setShortCitation (citation)
setTitle (title)
setType (reference_type)
setYear (year)
```

dipper.sources package

Submodules

dipper.sources.AnimalQTLdb module

```
class dipper.sources.AnimalQTLdb (graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source.Source
```

The Animal Quantitative Trait Loci (QTL) database (Animal QTLdb) is designed to house publicly all available QTL and single-nucleotide polymorphism/gene association data on livestock animal species. This includes:

- chicken
- horse
- cow
- sheep
- rainbow trout
- pig

While most of the phenotypes here are related to animal husbandry, production, and rearing, integration of these phenotypes with other species may lead to insight for human disease.

Here, we use the QTL genetic maps and their computed genomic locations to create associations between the QTLs and their traits. The traits come in their internal Animal Trait ontology vocabulary, which they further map to [Vertebrate Trait](<http://bioportal.bioontology.org/ontologies/VT>), Product Trait, and Clinical Measurement Ontology vocabularies.

Since these are only associations to broad locations, we link the traits via “is_marker_for”, since there is no specific causative nature in the association. p-values for the associations are attached to the Association objects. We default to the UCSC build for the genomic coordinates, and make equivalences.

Any genetic position ranges that are <0, we do not include here.

```
fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
        None

files = {'cattle_bp': {'file': 'QTL_Btau_4.6.gff.txt.gz', 'url': 'http://www.animal}

getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse(limit=None)
    Parameters limit -
    Returns

test_ids = {1795, 28483, 32133, 1798, 29385, 29018, 31023, 8945, 17138, 12532, 29016,
```

dipper.sources.Bgee module

```

class dipper.sources.Bgee.Bgee(graph_type, are_bnodes_skolemized, tax_ids=None, version=None)
Bases: dipper.sources.Source.Source

Bgee is a database to retrieve and compare gene expression patterns between animal species.

Bgee first maps heterogeneous expression data (currently RNA-Seq, Affymetrix, in situ hybridization, and EST data) to anatomy and development of different species.

Then, in order to perform automated cross species comparisons, homology relationships across anatomies, and comparison criteria between developmental stages, are designed.

BGEE_FTP = 'ftp.bgee.org'
DEFAULT_TAXA = [10090, 10116, 13616, 28377, 6239, 7227, 7955, 8364, 9031, 9258, 9544,
checkIfRemoteIsNewer(localfile, remote_size, remote_modify)
    Overrides checkIfRemoteIsNewer in Source class

Parameters

- localfile – str file path
- remote_size – str bytes
- remote_modify – str last modify date in the form 20160705042714

Returns boolean True if remote file is newer else False

fetch(is_dl_forced=False)
    Parameters is_dl_forced – boolean, force download

    Returns

files = {'anat_entity': {'pattern': re.compile('.*_all_data_.*'), 'path': '/downloaded'}
parse(limit=None)
    Given the input taxa, expects files in the raw directory with the name
    {tax_id}_anat_entity_all_data_Pan_troglodytes.tsv.zip

    Parameters limit – int Limit to top ranked anatomy associations per group

    Returns None

```

dipper.sources.BioGrid module

```
class dipper.sources.BioGrid.BioGrid(graph_type, are_bnodes_skolemized, tax_ids=None)
    Bases: dipper.sources.Source.Source

    Biogrid interaction data

    biogrid_ids = [106638, 107308, 107506, 107674, 107675, 108277, 108506, 108767, 108814,
    fetch(is_dl_forced=False)

        Parameters is_dl_forced -

        Returns None

    files = {'identifiers': {'file': 'identifiers.tab.zip', 'url': 'http://thebiogrid.o
    getTestSuite()
        An abstract method that should be overwritten with tests appropriate for the specific source. :return:

    parse(limit=None)

        Parameters limit -

        Returns
```

dipper.sources.CTD module

```
class dipper.sources.CTD.CTD(graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.Source.Source
```

The Comparative Toxicogenomics Database (CTD) includes curated data describing cross-species chemical–gene/protein interactions and chemical– and gene–disease associations to illuminate molecular mechanisms underlying variable susceptibility and environmentally influenced diseases.

Here, we fetch, parse, and convert data from CTD into triples, leveraging only the associations based on DIRECT evidence (not using the inferred associations). We currently process the following associations: * chemical-disease * gene-pathway * gene-disease

CTD curates relationships between genes and chemicals/diseases with marker/mechanism and/or therapeutic. Unfortunately, we cannot disambiguate between marker (gene expression) and mechanism (causation) for these associations. Therefore, we are left to relate these simply by “marker”.

CTD also pulls in genes and pathway membership from KEGG and REACTOME. We create groups of these following the pattern that the specific pathway is a subclass of ‘cellular process’ (a go process), and the gene is “involved in” that process.

For diseases, we preferentially use OMIM identifiers when they can be used uniquely over MESH. Otherwise, we use MESH ids.

Note that we scrub the following identifiers and their associated data: * REACT:REACT_116125 - generic disease class * MESH:D004283 - dog diseases * MESH:D004195 - disease models, animal * MESH:D030342 - genetic diseases, inborn * MESH:D040181 - genetic diseases, x-linked * MESH:D020022 - genetic predisposition to a disease

```
fetch(is_dl_forced=False)
```

Override Source.fetch() Fetches resources from CTD using the CTD.files dictionary Args: :param is_dl_forced (bool): Force download Returns: :return None

```
files = {'chemical_disease_interactions': {'file': 'CTD_chemicals_diseases.tsv.gz',
```

```
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (limit=None)
    Override Source.parse() Parses version and interaction information from CTD Args: :param limit (int, optional) limit the number of rows processed Returns: :return None

static_files = {'publications': {'file': 'CTD_curated_references.tsv'}}
```

dipper.sources.ClinVar module

```
class dipper.sources.ClinVar(graph_type, are_bnodes_skolemized, tax_ids=None, gene_ids=None)
Bases: dipper.sources.Source

ClinVar is a host of clinically relevant variants, both directly-submitted and curated from the literature. We process the variant_summary file here, which is a digested version of their full xml. We add all variants (and coordinates/build) from their system.

fetch (is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
    None

files = {'variant_citations': {'file': 'variant_citations.txt', 'url': 'http://ftp..}}
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

scrub ()
    The var_citations file has a bad row in it with > 6 cols. I will comment these out.

Returns

variant_ids = [4288, 4289, 4290, 4291, 4297, 5240, 5241, 5242, 5243, 5244, 5245, 5246,
```

dipper.sources.ClinVarXML_alpha module

dipper.sources.Coriell module

```
class dipper.sources.Coriell(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source
```

The Coriell Catalog provided to Monarch includes metadata and descriptions of NIGMS, NINDS, NHGRI, and NIA cell lines. These lines are made available for research purposes. Here, we create annotations for the cell lines as models of the diseases from which they originate.

We create a handle for a patient from which the given cell line is derived (since there may be multiple cell lines created from a given patient). A genotype is assembled for a patient, which includes a karyotype (if specified) and/or a collection of variants. Both the genotype (has_genotype) and disease are linked to the patient (has_phenotype), and the cell line is listed as derived from the patient. The cell line is classified by its [CLO cell type](<http://www.ontobee.org/browser/index.php?o=clo>), which itself is linked to a tissue of origin.

Unfortunately, the omim numbers listed in this file are both for genes & diseases; we have no way of knowing a priori if a designated omim number is a gene or disease; so we presently link the patient to any omim id via the has_phenotype relationship.

Notice: The Coriell catalog is delivered to Monarch in a specific format, and requires ssh rsa fingerprint identification. Other groups wishing to get this data in it's raw form will need to contact Coriell for credential. This needs to be placed into your configuration file for it to work.

fetch (is_dl_forced=False)

Here we connect to the coriell sftp server using private connection details. They dump bi-weekly files with a timestamp in the filename. For each catalog, we poll the remote site and pull the most-recently updated file, renaming it to our local latest.csv.

Be sure to have pg user/password connection details in your conf.json file, like: dbauth : {“coriell” : {“user” : “<username>”, “password” : “<password>”, “host” : <host>, “private_key”=path/to/rsa_key} }

Parameters is_dl_forced –

Returns

```
files = {'NHGRI': {'file': 'NHGRI.csv', 'page': 'https://catalog.coriell.org/1/NHGRI'}}

getTestSuite()
An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (limit=None)
abstract method to parse all data from an external resource, that was fetched in fetch() this should be
overridden by subclasses :return: None

terms = {'age': 'EFO:0000246', 'cell_line_repository': 'CLO:0000008', 'collection':
test_lines = ['ND02380', 'ND02381', 'ND02383', 'ND02384', 'GM17897', 'GM17898', 'GM17899']
```

dipper.sources.Decipher module

```
class dipper.sources.Decipher(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source.Source
```

The Decipher group curates and assembles the Development Disorder Genotype Phenotype Database (DDG2P) which is a curated list of genes reported to be associated with developmental disorders, compiled by clinicians as part of the DDD study to facilitate clinical feedback of likely causal variants.

Beware that the redistribution of this data is a bit unclear from the [license](<https://decipher.sanger.ac.uk/legal>). If you intend to distribute this data, be sure to have the appropriate licenses in place.

fetch (is_dl_forced=False)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
None

```
files = {'annot': {'file': 'ddg2p.zip', 'url': 'https://decipher.sanger.ac.uk/files'}}
```

make_allele_by_consequence (consequence, gene_id, gene_symbol)

Given a “consequence” label that describes a variation type, create an anonymous variant of the specified gene as an instance of that consequence type.

Parameters

- **consequence** –
- **gene_id** –
- **gene_symbol** –

Returns allele_id

parse (*limit=None*)
abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

dipper.sources.EOM module

```
class dipper.sources.EOM(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.PostgreSQLSource.PostgreSQLSource
```

Elements of Morphology is a resource from NHGRI that has definitions of morphological abnormalities, together with image depictions. We pull those relationships, as well as our local mapping of equivalences between EOM and HP terminologies.

The website is crawled monthly by NIF's DISCO crawler system, which we utilize here. Be sure to have pg user/password connection details in your conf.json file, like: dbauth : { 'disco' : { 'user' : '<username>', 'password' : '<password>' } }

Monarch-curated data for the HP to EOM mapping is stored at <https://raw.githubusercontent.com/obophenotype/human-phenotype-ontology/master/src/mappings/hp-to-eom-mapping.tsv>

Since this resource is so small, the entirety of it is the “test” set.

```
fetch(is_dl_forced=False)
    create the connection details for DISCO
files = {'map': {'file': 'hp-to-eom-mapping.tsv', 'url': 'https://raw.githubusercontent.com/obophenotype/human-phenotype-ontology/master/src/mappings/hp-to-eom-mapping.tsv'}}
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:
parse(limit=None)
    Over ride Source.parse inherited via PostgreSQLSource
tables = ['dvp.pr_nlx_157874_1']
```

dipper.sources.Ensembl module

```
class dipper.sources.Ensembl.Ensembl(graph_type, are_bnodes_skolemized, tax_ids=None, gene_ids=None)
Bases: dipper.sources.Source.Source
```

This is the processing module for Ensembl.

It only includes methods to acquire the equivalences between NCBIGene and ENSG ids using ENSEMBL's Biomart services.

```
fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None
fetch_protein_gene_map(taxon_id)
    Fetch a list of proteins for a species in biomart :param taxid: :return: dict
fetch_protein_list(taxon_id)
    Fetch a list of proteins for a species in biomart :param taxid: :return: list
fetch_uniprot_gene_map(taxon_id)
    Fetch a dict of uniprot-gene for a species in biomart :param taxid: :return: dict
```

```
files = {'10090': {'file': 'ensembl_10090.txt'}, '10116': {'file': 'ensembl_10116.txt'}}

getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return: TestSuite

parse(limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be
    overridden by subclasses :return: None
```

dipper.sources.FlyBase module

```
class dipper.sources.FlyBase.FlyBase(graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.PostgreSQLSource.PostgreSQLSource
```

This is the [Drosophila Genetics](<http://www.flybase.org/>) resource, from which we process genotype and phenotype data about fruitfly. Genotypes leverage the GENO genotype model.

Here, we connect to their public database, and download a subset of tables/views to get specifically at the geno-pheno data, then iterate over the tables. We end up effectively performing joins when adding nodes to the graph. We connect using the [Direct Chado Access](http://gmod.org/wiki/Public_Chado_Databases#Direct_Chado_Access)

When running the whole set, it performs best by dumping raw triples using the flag `--format nt``.

fetch (*is_dl_forced=False*)

Returns

```
files = {'disease_models': { 'file': 'allele_human_disease_model_data.tsv.gz', 'url':
```

`getTestSuite()`

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (*limit=None*)

We process each of the postgres tables in turn. The order of processing is important here, as we build up a hashmap of internal vs external identifiers (unique keys by type to FB id). These include allele, marker (gene), publication, strain, genotype, annotation (association), and descriptive notes. :param limit: Only parse this many lines of each table :return:

```
queries = {'feature': "\n    SELECT feature_id, dbxref_id, organism_id, name, uniquename,\nresources = [{query': '../resources/sql/fb/feature_relationship.sql', 'outfile':\ntables = ['genotype', 'feature_genotype', 'pub', 'feature_pub', 'pub_dbxref', 'feature'\ntest_keys = {'allele': [29677937, 23174110, 23230960, 23123654, 23124718, 23146222, 2
```

dipper.sources.GWASCatalog module

```
class dipper.sources.GWASCatalog.GWASCatalog(graph_type, are_bnodes_skolemized)  
    Bases: dipper.sources.Source.Source
```

The NHGRI-EBI Catalog of published genome-wide association studies.

We link the variants recorded here to the curated EFO-classes using a “contributes_to” linkage because the only thing we know is that the SNPs are associated with the trait/disease, but we don’t know if it is actually causative.

Description of the GWAS catalog is here: http://www.ebi.ac.uk/gwas/docs/fileheaders#_file_headers_for_catalog_version_1_0_1

GWAS also publishes Owl files described here <http://www.ebi.ac.uk/gwas/docs/ontology>

Status: IN PROGRESS

```
GWASFILE = 'gwas-catalog-associations_ontology-annotated.tsv'
GWASFTP = 'ftp://ftp.ebi.ac.uk/pub/databases/gwas/releases/latest'
fetch(is_dl_forced=False)
    Parameters is_dl_forced –
    Returns
files = {'catalog': {'file': 'gwas-catalog-associations_ontology-annotated.tsv', 'url': 'http://ftp.ebi.ac.uk/pub/databases/gwas/releases/latest/gwas-catalog-associations_ontology-annotated.tsv'}}
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:
parse(limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None
process_catalog(limit=None)
    Parameters limit –
    Returns
terms = {'age': 'EFO:0000246', 'cell_line_repository': 'CLO:0000008', 'collection': 'GO:0000001'}
```

dipper.sources.GeneOntology module

```
class dipper.sources.GeneOntology(graph_type, are_bnodes_skolemized, tax_ids=None)
Bases: dipper.sources.Source
```

This is the parser for the [Gene Ontology Annotations](<http://www.geneontology.org>), from which we process gene-process/function/subcellular location associations.

We generate the GO graph to include the following information: * genes * gene-process * gene-function * gene-location

We process only a subset of the organisms:

Status: IN PROGRESS / INCOMPLETE

```
clean_db_prefix(db)
```

Here, we map the GO-style prefixes with Monarch-style prefixes that are able to be processed by our curie_map. :param db: :return:

```
fetch(is_dl_forced=False)
```

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

```
files = {'10090': {'file': 'gene_association.mgi.gz', 'url': 'http://geneontology.org/gene_association.mgi.gz'}}
```

```
getTestSuite()
```

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

```
get_uniprot_entrez_id_map()
```

```
map_files = {'eco_map': 'http://purl.obolibrary.org/obo/eco/gaf-eco-mapping.txt'}
```

```
parse(limit=None)
```

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

process_gaf (*file, limit, id_map=None, eco_map=None*)

dipper.sources.GeneReviews module

class `dipper.sources.GeneReviews` (*graph_type, are_bnodes_skolemized*)
Bases: `dipper.sources.Source`

Here we process the GeneReviews mappings to OMIM, plus inspect the GeneReviews (html) books to pull the clinical descriptions in order to populate the definitions of the terms in the ontology. We define the GeneReviews items as classes that are either grouping classes over OMIM disease ids (gene ids are filtered out), or are made as subclasses of DOID:4 (generic disease).

Note that GeneReviews [copyright policy](<http://www.ncbi.nlm.nih.gov/books/NBK138602/>) (as of 2015.11.20) says:

GeneReviews® chapters are owned by the University of Washington, Seattle, © 1993-2015. Permission is hereby granted to reproduce, distribute, and translate copies of content materials provided that (i) credit for source (www.ncbi.nlm.nih.gov/books/NBK1116/) and copyright (University of Washington, Seattle) are included with each copy; (ii) a link to the original material is provided whenever the material is published elsewhere on the Web; and (iii) reproducers, distributors, and/or translators comply with this copyright notice and the GeneReviews Usage Disclaimer.

This script doesn't pull the GeneReviews books from the NCBI Bookshelf directly; scripting this task is expressly prohibited by [NCBIBookshelf policy](<http://www.ncbi.nlm.nih.gov/books/NBK45311/>). However, assuming you have acquired the books (in html format) via permissible means, a parser for those books is provided here to extract the clinical descriptions to define the NBK identified classes.

create_books ()

fetch (*is_dl_forced=False*)

We fetch GeneReviews id-label map and id-omim mapping files from NCBI. :return: None

files = {'idmap': {'file': 'NBKid_shortname_OMIM.txt', 'url': '<http://ftp.ncbi.nih.gov/genereviews/books/>'}}

getTestSuite ()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (*limit=None*)

Returns None

process_nbk_html (*limit*)

Here we process the gene reviews books to fetch the clinical descriptions to include in the ontology.

We only use books that have been acquired manually, as NCBI Bookshelf does not permit automated downloads. This parser will only process the books that are found in the `raw/genereviews/books` directory, permitting partial completion.

Parameters *limit* –

Returns

dipper.sources.HGNC module

class `dipper.sources.HGNC` (*graph_type, are_bnodes_skolemized, tax_ids=None, gene_ids=None*)
Bases: `dipper.sources.Source`

This is the processing module for HGNC.

We create equivalences between HGNC identifiers and ENSEMBL and NCBI Gene. We also add the links to cytogenic locations for the gene features.

```
fetch (is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
        None

files = {'genes': {'file': 'hgnc_complete_set.txt', 'url': 'ftp://ftp.ebi.ac.uk/pub/...'}
getTestSuite ()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

get_symbol_id_map ()
    A convenience method to create a mapping between the HGNC symbols and their identifiers. :return:

parse (limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be
    overridden by subclasses :return: None
```

dipper.sources.HPOAnnotations module

```
class dipper.sources.HPOAnnotations (graph_type,
                                    are_bnodes_skolemized)
Bases: dipper.sources.Source.Source
```

The [Human Phenotype Ontology](<http://human-phenotype-ontology.org>) group curates and assembles over 115,000 annotations to hereditary diseases using the HPO ontology. Here we create OBAN-style associations between diseases and phenotypic features, together with their evidence, and age of onset and frequency (if known). The parser currently only processes the “abnormal” annotations. Association to “remarkable normality” will be added in the near future.

We create additional associations from text mining. See info at <http://pubmed-browser.human-phenotype-ontology.org/>.

Also, you can read about these annotations in [PMID:26119816](<http://www.ncbi.nlm.nih.gov/pubmed/26119816>).

In order to properly test this class, you should have a conf.json file configured with some test ids, in the structure of: # as examples. put your favorite ids in the config. <pre> test_ids: {"disease": ["OMIM:119600", "OMIM:120160"]} </pre>

```
add_common_files_to_file_list ()
eco_dict = {'ICE': 'ECO:0000305', 'IEA': 'ECO:0000501', 'ITM': 'ECO:0000246', 'PCS': '...'}
fetch (is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
        None

files = {'annot': {'file': 'phenotype_annotation.tab', 'url': 'http://compbio.chari...'}
getTestSuite ()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

get_common_files ()
    Fetch the raw hpo-annotation-data by cloning/pulling the [repository](https://github.com/monarch-initiative/hpo-annotation-data.git) These files get added to the files object, and iterated
    over separately. :return:

get_doid_ids_for_unpadding ()
    Here, we fetch the doid owl file, and get all the doids. We figure out which are not zero-padded, so we can
    map the DOID to the correct identifier when processing the common annotation files.
```

This may become obsolete when <https://github.com/monarch-initiative/hpo-annotation-data/issues/84> is addressed.

Returns

`parse (limit=None)`

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

`process_all_common_disease_files (limit=None)`

Loop through all of the files that we previously fetched from git, creating the disease-phenotype assoc. :param limit: :return:

`process_common_disease_file (raw, unpadded_doids, limit=None)`

Make disease-phenotype associations. Some identifiers need clean up: * DOIDs are listed as DOID-DOID: -> DOID: * DOIDs may be unnecessarily zero-padded. these are remapped to their non-padded equivalent.

Parameters

- `raw` –
- `unpadded_doids` –
- `limit` –

Returns

`scrub ()`

Perform various data-scrubbing on the raw data files prior to parsing. For this resource, this currently includes: * revise errors in identifiers for some OMIM and PMIDs

Returns

dipper.sources.IMPC module

`class dipper.sources.IMPC(graph_type, are_bnodes_skolemized)`

Bases: [dipper.sources.Source](#)

From the [IMPC](<http://mousephenotype.org>) website: The IMPC is generating a knockout mouse strain for every protein coding gene by using the embryonic stem cell resource generated by the International Knockout Mouse Consortium (IKMC). Systematic broad-based phenotyping is performed by each IMPC center using standardized procedures found within the International Mouse Phenotyping Resource of Standardised Screens (IMPRess) resource. Gene-to-phenotype associations are made by a versioned statistical analysis with all data freely available by this web portal and by several data download features.

Here, we pull the data and model the genotypes using GENO and the genotype-to-phenotype associations using the OBAN schema.

We use all identifiers given by the IMPC with a few exceptions:

- For identifiers that IMPC provides, but does not resolve,

we instantiate them as Blank Nodes. Examples include things with the pattern of: UROALL, EUROCURATE, NULL-* ,

- We mint three identifiers:

1. Intrinsic genotypes not including sex, based on:

- colony_id (ES cell line + phenotyping center)

- strain
 - zygosity

2. For the Effective genotypes that are attached to the phenotypes:

- colony_id (ES cell line + phenotyping center)
 - strain
 - zygosity
 - sex

3. Associations based on: effective_genotype_id + phenotype_id + phenotyping_center + pipeline_stable_id + procedure_stable_id + parameter_stable_id

We DO NOT yet add the assays as evidence for the G2P associations here. To be added in the future.

compare_checksums()

test to see if fetched file matches checksum from ebi :return: True or False

fetch (*is_dl_forced=False*)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

```
files = {'all': {'file': 'ALL_genotype_phenotype.csv.gz', 'url': 'ftp://ftp.ebi.ac.uk/pub/databases/ukb/phenotype/ALL_genotype_phenotype.csv.gz'}}
```

getTestSuite()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

```
map_files = {'impc map': '.../resources/impc_mappings.yaml', 'impress map': 'https://raw.githubusercontent.com/AllenInstitute/impress/master/mappings.yaml'}
```

parse (*limit=None*)

IMPC data is delivered in three separate csv files OR in one integrated file, each with the same file format.

Parameters `limit` –

Returns

Returns

:param file :return dict

.param file .return dict

test_ids = [MG1.109380 , MG1.1347004 , MG1.1333493 , MG1.1913840 , MG1.2144137 ,

dipper.sources.KEGG module

class dipper.sources.KEGG.KEGG(*graph_type*, *are_bnodes_skolemized*)
 Base class for KEGG sources.

Bases: `dipper.sources.Source.Source`

```
fetch(is_dl_forced=False)
```

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

```
files = {'cel_orthologs': {'file': 'cel_orthologs', 'url': 'http://rest.kegg.jp/lin'}}
```

getTestSuite()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (*limit=None*)

Parameters limit –

Returns

```
test_ids = {'disease': ['ds:H00015', 'ds:H00026', 'ds:H00712', 'ds:H00736', 'ds:H00011']}
```

dipper.sources.MGI module

```
class dipper.sources.MGI.MGI(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.PostgreSQLSource.PostgreSQLSource
```

This is the [Mouse Genome Informatics](<http://www.informatics.jax.org/>) resource, from which we process genotype and phenotype data about laboratory mice. Genotypes leverage the GENO genotype model.

Here, we connect to their public database, and download a subset of tables/views to get specifically at the geno-pheno data, then iterate over the tables. We end up effectively performing joins when adding nodes to the graph. In order to use this parser, you will need to have user/password connection details in your conf.json file, like: dbauth : {‘mgi’ : {‘user’ : ‘<username>’, ‘password’ : ‘<password>’}} You can request access by contacting mgi-help@jax.org

```
fetch(is_dl_forced=False)
```

For the MGI resource, we connect to the remote database, and pull the tables into local files. We’ll check the local table versions against the remote version :return:

```
fetch_transgene_genes_from_db(cxn)
```

This is a custom query to fetch the non-mouse genes that are part of transgene alleles.

Parameters cxn –

Returns

```
getTestSuite()
```

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

```
parse(limit=None)
```

We process each of the postgres tables in turn. The order of processing is important here, as we build up a hashmap of internal vs external identifiers (unique keys by type to MGI id). These include allele, marker (gene), publication, strain, genotype, annotation (association), and descriptive notes. :param limit: Only parse this many lines of each table :return:

```
process_mgi_note_allele_view(limit=None)
```

These are the descriptive notes about the alleles. Note that these notes have embedded HTML - should we do anything about that? :param limit: :return:

```
process_mgi_relationship_transgene_genes(limit=None)
```

Here, we have the relationship between MGI transgene alleles, and the non-mouse gene ids that are part of them. We augment the allele with the transgene parts.

Parameters limit –

Returns

```
resources = [{‘query’: ‘.../..../resources/sql/mgi/mgi_dbinfo.sql’, ‘outfile’: ‘mgi_dbinfo.out’}]
```

```
test_keys = {‘allele’: [1612, 1609, 1303, 56760, 816699, 51074, 14595, 816707, 246, 300]}
```

dipper.sources.MGISlim module

```
class dipper.sources.MGISlim.MGISlim(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source.Source
```

slim mgf model only containing Gene to phenotype associations Uses mousemine: <http://www.mousemine.org/mousemine/begin.do>

fetch (is_dl_forced)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

parse (limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

dipper.sources.MMRRC module

class dipper.sources.MMRRC (*graph_type*, *are_bnodes_skolemized*)

Bases: *dipper.sources.Source*

Here we process the Mutant Mouse Resource and Research Center (<https://www.mmrrc.org>) strain data, which includes: * strains, their mutant alleles * phenotypes of the alleles * descriptions of the research uses of the strains

Note that some gene identifiers are not included (for many of the transgenics with human genes) in the raw data. We do our best to process the links between the variant and the affected gene, but sometimes the mapping is not clear, and we do not include it. Many of these details will be solved by merging this source with the MGI data source, who has the variant-to-gene designations.

Also note that even though the strain pages at the MMRC site do list phenotypic differences in the context of the strain backgrounds, they do not provide that data to us, and thus we cannot supply that disambiguation here.

fetch (is_dl_forced=False)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

files = {'catalog': {'file': 'mmrrc_catalog_data.csv', 'url': 'https://www.mmrrc.org'}}

getTestSuite ()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

parse (limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

test_ids = ['MMRRC:037507-MU', 'MMRRC:041175-UCD', 'MMRRC:036933-UNC', 'MMRRC:037884-UCD']

dipper.sources.MPD module

class dipper.sources.MPD (*graph_type*, *are_bnodes_skolemized*)

Bases: *dipper.sources.Source*

From the [MPD](<http://phenome.jax.org/>) website: This resource is a collaborative standardized collection of measured data on laboratory mouse strains and populations. Includes baseline phenotype data sets as well as studies of drug, diet, disease and aging effect. Also includes protocols, projects and publications, and SNP, variation and gene expression studies.

Here, we pull the data and model the genotypes using GENO and the genotype-to-phenotype associations using the OBAN schema.

MPD provide measurements for particular assays for several strains. Each of these measurements is itself mapped to a MP or VT term as a phenotype. Therefore, we can create a strain-to-phenotype association based on those strains that lie outside of the “normal” range for the given measurements. We can compute the average of the measurements for all strains tested, and then threshold any extreme measurements being beyond some threshold beyond the average.

Our default threshold here, is +/- standard deviations beyond the mean.

Because the measurements are made and recorded at the level of a specific sex of each strain, we associate the MP/VT phenotype with the sex-qualified genotype/strain.

```
MPDDL = 'http://phenomedoc.jax.org/MPD_downloads'

static build_measurement_description(row)

static check_header(filename, header)

fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
    None

files = {'assay_metadata': {'file': 'measurements.csv', 'url': 'http://phenomedoc.j
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

mgd_agent_id = 'MPD:db/q?rtn=people/allinv'
mgd_agent_label = 'Mouse Phenotype Database'
mgd_agent_type = 'foaf:organization'

static normalise_units(units)

parse(limit=None)
    MPD data is delivered in four separate csv files and one xml file, which we process iteratively and write
    out as one large graph.

    Parameters limit -
    Returns

test_ids = ['MPD:6', 'MPD:849', 'MPD:425', 'MPD:569', 'MPD:10', 'MPD:1002', 'MPD:39',
```

dipper.sources.Monarch module

```
class dipper.sources.Monarch(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source
```

This is the parser for data curated by the [Monarch Initiative](<https://monarchinitiative.org>). Data is currently maintained in a private repository, soon to be released.

```
fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
    None

parse(limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be
    overridden by subclasses :return: None

process_omia_phenotypes(limit)
```

dipper.sources.Monochrom module

```
class dipper.sources.Monochrom(graph_type, are_bnodes_skolemized,  
                                tax_ids=None)  
Bases: dipper.sources.Source.Source
```

This class will leverage the GENO ontology and modeling patterns to build an ontology of chromosomes for any species. These classes represent major structural pieces of Chromosomes which are often universally referenced, using physical properties/observations that remain constant over different genome builds (such as banding patterns and arms). The idea is to create a scaffold upon which we can hang build-specific chromosomal coordinates, and reason across them.

In general, this will take the cytogenic bands files from UCSC, and create missing grouping classes, in order to build the partonomy from a very specific chromosomal band up through the chromosome itself and enable overlap and containment queries. We use RO:subsequence_of as our relationship between nested chromosomal parts. For example, 13q21.31 ==> 13q21.31, 13q21.3, 13q21, 13q2, 13q, 13

At the moment, this only computes the bands for Human, Mouse, Zebrafish, and Rat but will be expanding in the future as needed.

Because this is a universal framework to represent the chromosomal structure of any species, we must mint identifiers for each chromosome and part. We differentiate species by first creating a species-specific genome, then for each species-specific chromosome we include the NCBI taxon number together with the chromosome number, like: `<species number>chr<num><band>`. For 13q21.31, this would be 9606chr13q21.31. We then create triples for a given band like: <pre> CHR:9606chr1p36.33 rdf[type] SO:chromosome_band CHR:9606chr1p36 subsequence_of :9606chr1p36.3 </pre> where any band in the file is an instance of a chr_band (or a more specific type), is a subsequence of it's containing region.

We determine the containing regions of the band by parsing the band-string; since each alphanumeric is a significant “place”, we can split it with the shorter strings being parents of the longer string

Since this is small, and we have not limited other items in our test set to a small region, we simply use the whole graph (genome) for testing purposes, and copy the main graph to the test graph.

Since this Dipper class is building an ONTOLOGY, rather than instance-level data, we must also include domain and range constraints, and other owl-isms.

TODO: any species by commandline argument

We are currently mapping these to the **CHR idspace**, but this is NOT YET APPROVED and is subject to change.

```
fetch(is_dl_forced=False)  
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:  
    None  
  
files = {'10090': {'file': '10090cytoBand.txt.gz', 'genome_label': 'Mouse', 'url':  
    getTestSuite()  
        An abstract method that should be overwritten with tests appropriate for the specific source. :return:  
    make_parent_bands(band, child_bands)  
        this will determine the grouping bands that it belongs to, recursively 13q21.31 ==> 13, 13q, 13q2, 13q21,  
        13q21.3, 13q21.31
```

Parameters

- **band** –
- **child_bands** –

Returns

map_type_of_region(regiontype)

Note that “stalk” refers to the short arm of acrocentric chromosomes chr13,14,15,21,22 for human. :param regiontype: :return:

parse(limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

region_type_map = {'acen': 'SO:0000577', 'chromosome': 'SO:0000340', 'chromosome_arm':

dipper.sources.Monochrom.**getChrPartTypeByNotation**(notation)

This method will figure out the kind of feature that a given band is based on pattern matching to standard karyotype notation. (e.g. 13q22.2 ==> chromosome sub-band)

This has been validated against human, mouse, fish, and rat nomenclature. :param notation: the band (without the chromosome prefix) :return:

dipper.sources.MyChem module

class dipper.sources.MyChem(graph_type, are_bnodes_skolemized)

Bases: *dipper.sources.Source*

static add_relation(results, relation)

static check_uniprot(target_dict)

static chunks(l, n)

Yield successive n-sized chunks from l.

static execute_query(query)

fetch(is_dl_forced=False)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

fetch_from_mychem()

static format_actions(target_dict)

static get_drug_record(ids, fields)

static get_inchikeys()

make_triples(source, package)

parse(limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

static return_target_list(targ_in)

dipper.sources.MyDrug module

class dipper.sources.MyDrug(graph_type, are_bnodes_skolemized)

Bases: *dipper.sources.Source*

Drugs and Compounds stored in the BioThings database

MY_DRUG_API = 'http://c.biotothings.io/v1/query'

checkIfRemoteIsNewer (*localfile*)

Need to figure out how biothings records releases, for now if the file exists we will assume it is a fully downloaded cache :param localfile: str file path :return: boolean True if remote file is newer else False

fetch (*is_dl_forced=False*)

Note there is a unpublished mydrug client that works like this: from mydrug import MyDrugInfo md = MyDrugInfo() r = list(md.query('_exists_:_aeolus', fetch_all=True))

Parameters **is_dl_forced** – boolean, force download

Returns

```
files = {'aeolus': {'file': 'aeolus.json'}}
```

parse (*limit=None, or_limit=1*)

Parse mydrug files :param limit: int limit json docs processed :param or_limit: int odds ratio limit :return: None

dipper.sources.NCBIGene module

```
class dipper.sources.NCBIGene(graph_type, are_bnodes_skolemized, tax_ids=None, gene_ids=None)
```

Bases: *dipper.sources.Source*.*Source*

This is the processing module for the National Center for Biotechnology Information. It includes parsers for the gene_info (gene names, symbols, ids, equivalent ids), gene history (alt ids), and gene2pubmed publication references about a gene.

This creates Genes as classes, when they are properly typed as such. For those entries where it is an ‘unknown significance’, it is added simply as an instance of a sequence feature. It will add equivalentClasses for a subset of external identifiers, including: ENSEMBL, HGMD, MGI, ZFIN, and gene product links for HPRD. They are additionally located to their Chromosomal band (until we process actual genomic coords in a separate file).

We process the genes from the filtered taxa, starting with those configured by default (human, mouse, fish). This can be overridden in the calling script to include additional taxa, if desired. The gene ids in the conf.json will be used to subset the data when testing.

All entries in the gene_history file are added as deprecated classes, and linked to the current gene id, with “replaced_by” relationships.

Since we do not know much about the specific link in the gene2pubmed; we simply create a “mentions” relationship.

```
SCIGRAPH_BASE = 'https://scigraph-ontology-dev.monarchinitiative.org/scigraph/graph/'
```

add_orthologs_by_gene_group (*graph, gene_ids*)

This will get orthologies between human and other vertebrate genomes based on the gene_group annotation pipeline from NCBI. More information can be learned here: <http://www.ncbi.nlm.nih.gov/news/03-13-2014-gene-provides-orthologs-regions/> The method for associations is described in [PMCID:3882889](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3882889/) == [PMID:24063302](http://www.ncbi.nlm.nih.gov/pubmed/24063302/). Because these are only between human and vertebrate genomes, they will certainly miss out on very distant orthologies, and should not be considered complete.

We do not run this within the NCBI parser itself; rather it is a convenience function for others parsers to call.

Parameters

- **graph** –

- **gene_ids** – Gene ids to fetch the orthology

Returns

```
fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
        None

files = {'gene2pubmed': {'file': 'gene2pubmed.gz', 'url': 'http://ftp.ncbi.nih.gov/'}

getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:

static map_type_of_gene(sotype)
parse(limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be
    overridden by subclasses :return: None

resources = {'clique_leader': '../..../resources/clique_leader.yaml'}
```

dipper.sources.OMA module

```
class dipper.sources.OMA.OMA(graph_type, are_bnodes_skolemized, tax_ids=None)
    Bases: dipper.sources.OrthoXML.OrthoXML

    BENCHMARK_BASE = 'https://omabrowser.org/ReferenceProteomes'

    files = {'oma_hogs': {'file': 'OMA_GETHOGs-2_2017-04.orthoxml.gz', 'url': 'https://'}

getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:
```

dipper.sources.OMIA module

```
class dipper.sources.OMIA.OMIA(graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.Source.Source
```

This is the parser for the [Online Mendelian Inheritance in Animals (OMIA)](<http://www.omnia.angis.org.au>), from which we process inherited disorders, other (single-locus) traits, and genes in >200 animal species (other than human and mouse and rats).

We generate the omia graph to include the following information: * genes * animal taxonomy, and breeds as instances of those taxa

(breeds are akin to “strains” in other taxa)

- animal diseases, along with species-specific subtypes of those diseases
- publications (and their mapping to PMIDs, if available)
- gene-to-phenotype associations (via an anonymous variant-locus
- breed-to-phenotype associations

We make links between OMIA and OMIM in two ways: 1. mappings between OMIA and OMIM are created as OMIA → hasdbXref OMIM 2. mappings between a breed and OMIA disease are created

to be a model for the mapped OMIM disease, IF AND ONLY IF it is a 1:1 mapping. there are some 1:many mappings, and these often happen if the OMIM item is a gene.

Because many of these species are not covered in the PANTHER orthology datafiles, we also pull any orthology relationships from the gene_group files from NCBI.

```
clean_up_omim_genes()
fetch(is_dl_forced=False)
    Parameters is_dl_forced –
    Returns
files = {'data': {'file': 'omia.xml.gz', 'url': 'http://compldb.angis.org.au/dumps/'}
getTestSuite()
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:
make_breed_id(key)
map_omia_group_category_to_ontology_id(category_num)
    Using the category number in the OMIA_groups table, map them to a disease id. This may be superceeded by other MONDO methods.
    Platelet disorders will be more specific once https://github.com/obophenotype/human-disease-ontology/issues/46 is fulfilled.
        Parameters category_num –
        Returns
parse(limit=None)
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None
process_associations(limit)
    Loop through the xml file and process the article-breed, article-phene, breed-phene, phene-gene associations, and the external links to LIDA.
        Parameters limit –
        Returns
process_classes(limit)
    Loop through the xml file and process the articles, breed, genes, phenes, and phenotype-grouping classes. We add elements to the graph, and store the id-to-label in the label_hash dict, along with the internal key-to-external id in the id_hash dict. The latter are referenced in the association processing functions.
        Parameters limit –
        Returns
process_species(limit)
    Loop through the xml file and process the species. We add elements to the graph, and store the id-to-label in the label_hash dict. :param limit: :return:
scrub()
    The XML file seems to have mixed-encoding; we scrub out the control characters from the file for processing.
    i.e.?i omia.xml:1555328.28: PCDATA invalid Char value 2 <field name="journal">Bulletin et Memoires de la Societe Centrale de Medic
        Returns
write_molgen_report()
```

dipper.sources.OMIM module

```
class dipper.sources.OMIM(graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.Source
```

The only anonymously obtainable data from the ftp site is mim2gene. However, more detailed information is available via their API. So, we pull the omim identifiers from their ftp site, then query their API in batches of 20. Their prescribed rate limits have been mercurial

one per two seconds or four per second, in 2017 November all mention of api rate limits have vanished
(save 20 IDs per call if any include is used)

Note this ingest requires an api Key which is not stored in the repo, but in a separate conf.json file.

Processing this source serves two purposes: 1. the creation of the OMIM classes for merging into the disease ontology 2. add annotations such as disease-gene associations

When creating the disease classes, we pull from their REST-api id/label/definition information. Additionally we pull the Orphanet and UMLS mappings (to make equivalent ids). We also pull the phenotypic series annotations as grouping classes.

```
fetch(is_dl_forced=True)
```

Get the preconfigured static files. This DOES NOT fetch the individual records via REST... that is handled in the parsing function. (To be refactored.) over riding Source.fetch() calling Source.get_files() :param is_dl_forced: :return:

```
files = {'all': {'file': 'mim2gene.txt', 'clean': 'https://data.omim.org/downloads/'}}
```

```
getTestSuite()
```

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

```
parse(limit=None)
```

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

```
process_entries(omimids, transform, included_fields=None, graph=None, limit=None)
```

Given a list of omim ids, this will use the omim API to fetch the entries, according to the `included_fields` passed as a parameter. If a transformation function is supplied, this will iterate over each entry, and either add the results to the supplied `graph` or will return a set of processed entries that the calling function can further iterate.

If no `included_fields` are provided, this will simply fetch the basic entry from omim, which includes an entry's: prefix, mimNumber, status, and titles.

Parameters

- **omimids** – the set of omim entry ids to fetch using their API
- **transform** – Function to transform each omim entry when looping
- **included_fields** – A set of what fields are required to retrieve from the API
- **graph** – the graph to add the transformed data into

Returns

```
test_ids = [119600, 120160, 157140, 158900, 166220, 168600, 219700, 253250, 305900, 60]
```

```
dipper.sources.OMIM.filter_keep_phenotype_entry_ids(entry, graph=None)
```

```
dipper.sources.OMIM.get_omim_id_from_entry(entry)
```

dipper.sources.Orphanet module

```
class dipper.sources.Orphanet.Orphanet (graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.Source
```

Orphanet's aim is to help improve the diagnosis, care and treatment of patients with rare diseases. For Orphanet, we are currently only parsing the disease-gene associations.

Note that ???

```
fetch (is_dl_forced=False)
```

Parameters **is_dl_forced** –

Returns

```
files = {'disease-gene': {'file': 'en_product6.xml', 'url': 'http://www.orphadata.o'}
```

```
getTestSuite()
```

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

```
parse (limit=None)
```

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

dipper.sources.OrthoXML module

```
class dipper.sources.OrthoXML.OrthoXML (graph_type, are_bnodes_skolemized, method,
                                         tax_ids=None)
    Bases: dipper.sources.Source
```

Extract the induced pairwise relations from an OrthoXML file.

This base class is primarily intended to extract the orthologous and paralogous relations from a file in OrthoXML file containing the QfO reference species data set.

A concrete method should subclass this class and overwrite the constructor method to provide the information about the dataset and a method name.

```
add_protein_to_graph
```

adds protein nodes to the graph and adds a “in_taxon” triple.

for efficiency reasons, we cache which proteins we have already added using a least recently used cache.

```
clean_protein_id (protein_id)
```

makes sure protein_id is properly prefixed

```
extract_taxon_info (gene_node)
```

extract the ncbi taxon id from a gene_node

default implementation goes up to the species node in the xml and extracts the id from the attribute at that node.

```
fetch (is_dl_forced=False)
```

Returns None

```
files = {}
```

```
parse (limit=None)
```

Returns None

```
class dipper.sources.OrthoXML.OrthoXMLParser (xml)
Bases: object

default_node_list()
extract_pairwise_relations (node=None)
get_children (node)
is_internal_node (node)
is_leaf (node)
leaf_label (leaf)
```

dipper.sources.Panther module

```
class dipper.sources.Panther.Panther (graph_type, are_bnodes_skolemized, tax_ids=None)
Bases: dipper.sources.Source.Source
```

The pairwise orthology calls from Panther DB: <http://pantherdb.org/> encompass 22 species, from the Re-Genome and HCOP projects. Here, we map the orthology classes to RO homology relationships. This resource may be extended in the future with additional species.

This currently makes a graph of orthologous relationships between genes, with the assumption that gene metadata (labels, equivalent ids) are provided from other sources.

Gene families are nominally created from the orthology files, though these are incomplete with no hierarchical (subfamily) information. This will get updated from the HMM files in the future.

Note that there is a fair amount of identifier cleanup performed to align with our standard CURIE prefixes.

The test graph of data is output based on configured “protein” identifiers in conf.json.

By default, this will produce a file with ALL orthologous relationships. IF YOU WANT ONLY A SUBSET, YOU NEED TO PROVIDE A FILTER UPON CALLING THIS WITH THE TAXON IDS

```
PNTHDL = 'ftp://ftp.pantherdb.org/ortholog/current_release'
fetch (is_dl_forced=False)

>Returns None

files = {'hcop': {'file': 'Orthologs_HCOP.tar.gz', 'url': 'ftp://ftp.pantherdb.org/'}
getTestSuite()
An abstract method that should be overwritten with tests appropriate for the specific source. :return:
parse (limit=None)

>Returns None
```

dipper.sources.PostgreSQLSource module

```
class dipper.sources.PostgreSQLSource.PostgreSQLSource (graph_type,
                                                       are_bnodes_skolemized,
                                                       name=None)
Bases: dipper.sources.Source.Source
```

Class for interfacing with remote Postgres databases

```
fetch_from_pgdb (tables, cxn, limit=None, force=False)
```

Will fetch all Postgres tables from the specified database in the cxn connection parameters.
This will save them to a local file named the same as the table, in tab-delimited format, including a header.

Parameters

- **tables** – Names of tables to fetch
- **cxn** – database connection details
- **limit** – A max row count to fetch for each table

Returns None

fetch_query_from_pgdb (*qname*, *query*, *con*, *cxn*, *limit=None*, *force=False*)

Supply either an already established connection, or connection parameters. The supplied connection will override any separate cxn parameter :param qname: The name of the query to save the output to :param query: The SQL query itself :param con: The already-established connection :param cxn: The postgres connection information :param limit: If you only want a subset of rows from the query :return:

dipper.sources.RGD module

class `dipper.sources.RGD(graph_type, are_bnodes_skolemized)`

Bases: `dipper.sources.Source`

Ingest of Rat Genome Database gene to mammalian phenotype gaf file

`RGD_BASE = 'ftp://ftp.rgd.mcw.edu/pub/data_release/annotated_rgd_objects_by_ontology/'`

fetch (*is_dl_forced=False*)

Override Source.fetch() Fetches resources from rat_genome_database using the rat_genome_database ftp site Args:

param is_dl_forced (bool) Force download

Returns: :return None

`files = {'rat_gene2mammalian_phenotype': {'file': 'rattus_genes_mp', 'url': 'ftp://.../rgd/gene2phenotype/rat_genes_mp'}}`

make_association (*record*)

construct the association :param record: :return: modeled association of genotype to mammalian phenotype

parse (*limit=None*)

Override Source.parse() Args:

:param limit (int, optional) limit the number of rows processed

Returns: :return None

dipper.sources.Reactome module

class `dipper.sources.Reactome.Reactome(graph_type, are_bnodes_skolemized)`

Bases: `dipper.sources.Source`

Reactome is a free, open-source, curated and peer reviewed pathway database. (<http://reactome.org/>)

`REACTOME_BASE = 'http://www.reactome.org/download/current/'`

```
fetch(is_dl_forced=False)
    Override Source.fetch() Fetches resources from reactome using the Reactome.files dictionary Args:
        param is_dl_forced (bool) Force download

    Returns: :return None

files = {'chebi2pathway': {'file': 'ChEBI2Reactome.txt', 'url': 'http://www.reactome.org/files/ChEBI2Reactome.txt'}
map_files = {'eco_map': 'http://purl.obolibrary.org/obo/eco/gaf-eco-mapping.txt'}
parse(limit=None)
    Override Source.parse() Args:
        :param limit (int, optional) limit the number of rows processed

    Returns: :return None
```

dipper.sources.SGD module

```
class dipper.sources.SGD(graph_type, are_bnodes_skolemized)
    Bases: dipper.sources.Source

    Ingest of Saccharomyces Genome Database (SGD) phenotype associations
    SGD_BASE = 'https://downloads.yeastgenome.org/curation/literature/'

    fetch(is_dl_forced=False)
        Override Source.fetch() Fetches resources from rat_genome_database using the rat_genome_database ftp site Args:
            param is_dl_forced (bool) Force download

    Returns: :return None

    files = {'sgd_phenotype': {'file': 'phenotype_data.tab', 'url': 'https://downloads.yeastgenome.org/curation/literature/phenotype_data.tab'}
    static make_apo_map()
    make_association(record)
        construct the association :param record: :return: modeled association of genotype to mammalian pheno-type
    parse(limit=None)
        Override Source.parse() Args:
            :param limit (int, optional) limit the number of rows processed

    Returns: :return None
```

dipper.sources.Source module

```
class dipper.sources.Source(graph_type, are_bnodes_skized=False, name=None)
    Bases: object

    Abstract class for any data sources that we'll import and process. Each of the subclasses will fetch() the data, scrub() it as necessary, then parse() it into a graph. The graph will then be written out to a single self.name().<dest_fmt> file.
```

checkIfRemoteIsNewer (*remote, local, headers*)

Given a remote file location, and the corresponding local file this will check the datetime stamp on the files to see if the remote one is newer. This is a convenience method to be used so that we don't have to re-fetch files that we already have saved locally :param *remote*: URL of file to fetch from remote server :param *local*: pathname to save file to locally :return: True if the remote file is newer and should be downloaded

compare_local_remote_bytes (*remotefile, localfile, remote_headers=None*)

test to see if fetched file is the same size as the remote file using information in the content-length field in the HTTP header :return: True or False

declareAsOntology (*graph*)

The file we output needs to be declared as an ontology, including it's version information.

TEC: I am not convinced dipper reformatting external data as RDF triples makes an OWL ontology (nor that it should be considered a goal).

Proper ontologies are built by ontologists. Dipper reforms data and annotates/decorates it with a minimal set of carefully arranged terms drawn from multiple proper ontologies. Which allows the whole (dipper's RDF triples and parent ontologies) to function as a single ontology we can reason over when combined in a store such as SciGraph.

Including more than the minimal ontological terms in dipper's RDF output constitutes a liability as it allows greater divergence between dipper artifacts and the proper ontologies.

Further information will be augmented in the dataset object. :param *version*: :return:

fetch (*is_dl_forced=False*)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

fetch_from_url (*remotefile, localfile=None, is_dl_forced=False, headers=None*)

Given a remote url and a local filename, attempt to determine if the remote file is newer; if it is, fetch the remote file and save it to the specified localfile, reporting the basic file information once it is downloaded :param *remotefile*: URL of remote file to fetch :param *localfile*: pathname of file to save locally :return: None

file_len (*fname*)**files** = {}**getTestSuite** ()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

static get_eco_map (*url*)

To convert the three column file to a hashmap we join primary and secondary keys, for example IEA GO_REF:0000002 ECO:0000256 IEA GO_REF:0000003 ECO:0000501 IEA Default ECO:0000501

becomes IEA-GO_REF:0000002: ECO:0000256 IEA-GO_REF:0000003: ECO:0000501 IEA: ECO:0000501

Returns dict

get_file_md5 (*directory, file, blocksize=1048576*)**get_files** (*is_dl_forced, files=None*)

Given a set of files for this source, it will go fetch them, and set a default version by date. If you need to set the version number by another method, then it can be set again. :param *is_dl_forced* - boolean :param *files* dict - override instance files dict :return: None

get_local_file_size (*localfile*)

Parameters **localfile** –

Returns size of file

get_remote_content_len(*remote*, *headers=None*)

Parameters *remote* –

Returns size of remote file

static hash_id(*long_string*)

prepend ‘b’ to avoid leading with digit truncate to 64bit sized word return sha1 hash of string :param long_string: str string to be hashed :return: str hash of id

static make_id(*long_string*, *prefix='MONARCH'*)

a method to create DETERMINISTIC identifiers based on a string’s digest. currently implemented with sha1 :param long_string: :return:

namespaces = {}

static open_and_parse_yaml(*file*)

Parameters *file* – String, path to file containing label-id mappings in the first two columns of each row

Returns dict where keys are labels and values are ids

parse(*limit*)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

static parse_mapping_file(*file*)

Parameters *file* – String, path to file containing label-id mappings in the first two columns of each row

Returns dict where keys are labels and values are ids

process_xml_table(*elem*, *table_name*, *processing_function*, *limit*)

This is a convenience function to process the elements of an xml document, when the xml is used as an alternative way of distributing sql-like tables. In this case, the “elem” is akin to an sql table, with it’s name of `table_name`. It will then process each `row` given the `processing_function` supplied.

Parameters

- **elem** – The element data
- **table_name** – The name of the table to process
- **processing_function** – The row processing function
- **limit** –

Appears to be making calls to the elementTree library although it not explicitly imported here.

Returns

static remove_backslash_r(*filename*, *encoding*)

A helpful utility to remove Carriage Return from any file. This will read a file into memory, and overwrite the contents of the original file.

TODO: This function may be a liability

Parameters *filename* –

Returns

```
settestmode (mode)
Set testMode to (mode). - True: run the Source in testMode; - False: run it in full mode :param mode: :return: None

settestonly (testonly)
Set that this source should only be processed in testMode :param testOnly: :return: None

whoami ()
write (fmt='turtle', stream=None)
This convenience method will write out all of the graphs associated with the source. Right now these are hardcoded to be a single "graph" and a "src_dataset.ttl" and a "src_test.ttl" If you do not supply stream='stdout' it will default write these to files.

In addition, if the version number isn't yet set in the dataset, it will be set to the date on file. :return: None
```

dipper.sources.StringDB module

```
class dipper.sources.StringDB.StringDB (graph_type, are_bnodes_skolemized, tax_ids=None,
                                         version=None)
Bases: dipper.sources.Source.Source
```

STRING is a database of known and predicted protein-protein interactions. The interactions include direct (physical) and indirect (functional) associations; they stem from computational prediction, from knowledge transfer between organisms, and from interactions aggregated from other (primary) databases. From: http://string-db.org/cgi/about.pl?footer_active_subpage=content

STRING uses one protein per gene. If there is more than one isoform per gene, we usually select the longest isoform, unless we have information that suggest that other isoform regarded as canonical (e.g., proteins in the CCDS database). From: <http://string-db.org/cgi/help.pl>

```
DEFAULT_TAXA = [9606, 10090, 7955, 7227, 6239]
```

```
STRING_BASE = 'http://string-db.org/download/'
```

```
fetch (is_dl_forced=False)
```

Override Source.fetch() Fetches resources from String

We also fetch ensembl to determine if protein pairs are from the same species Args:

```
    param is_dl_forced (bool) Force download
```

Returns: :return None

```
parse (limit=None)
```

Override Source.parse() Args:

```
    :param limit (int, optional) limit the number of rows processed
```

Returns: :return None

dipper.sources.UCSCBands module

```
class dipper.sources.UCSCBands.UCSCBands (graph_type, are_bnodes_skolemized,
                                              tax_ids=None)
Bases: dipper.sources.Source.Source
```

This will take the UCSC definitions of cytogenic bands and create the nested structures to enable overlap and containment queries. We use `Monochrom.py` to create the OWL-classes of the chromosomal parts. Here, we simply worry about the instance-level values for particular genome builds.

Given a chr band definition, the nested containment structures look like: 13q21.31 => 13q21.31, 13q21.3, 13q21, 13q2, 13q, 13

We determine the containing regions of the band by parsing the band-string; since each alphanumeric is a significant “place”, we can split it with the shorter strings being parents of the longer string. # Here we create build-specific chroms, which are instances of the classes produced from `Monochrom.py`. You can instantiate any number of builds for a genome.

We leverage the Faldo model here for region definitions, and map each of the chromosomal parts to SO.

We differentiate the build by adding the build id to the identifier prior to the chromosome number. These then are instances of the species-specific chromosomal class.

The build-specific chromosomes are created like: <pre> <build number>chr<num><band> with triples for a given band like: _:hg19chr1p36.33

```
rdfs:type SO:chromosome_band, faldo:Region, CHR:9606chr1p36.33, subsequence_of  
_:hg19chr1p36.3, faldo:location [ a faldo:BothStrandPosition  
faldo:begin 0, faldo:end 2300000, faldo:reference 'hg19'  
].
```

</pre> where any band in the file is an instance of a chr_band (or a more specific type), is a subsequence of it's containing region, and is located in the specified coordinates.

We do not have a separate graph for testing.

TODO: any species by commandline argument

```
HGGP = 'http://hgdownload.cse.ucsc.edu/goldenPath'  
fetch(is_dl_forced=False)  
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:  
    None  
  
files = {'10090': {'file': 'mm10cytoband.txt.gz', 'genome_label': 'Mouse', 'url':  
getTestSuite()  
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:  
parse(limit=None)  
    abstract method to parse all data from an external resource, that was fetched in fetch() this should be  
    overridden by subclasses :return: None
```

dipper.sources.UDP module

```
class dipper.sources.UDP(graph_type, are_bnodes_skolemized)  
Bases: dipper.sources.Source.Source
```

The National Institutes of Health (NIH) Undiagnosed Diseases Program (UDP) is part of the Undiagnosed Disease Network (UDN), an NIH Common Fund initiative that focuses on the most puzzling medical cases referred to the NIH Clinical Center in Bethesda, Maryland. from <https://www.genome.gov/27544402/the-undiagnosed-diseases-program/>

Data is available by request for access via the NHGRI collaboration server: <https://udplims-collab.nhgri.nih.gov/api>

Note the fetcher requires credentials for the UDP collaboration server. Credentials are added via a config file, config.json, in the following format {

```
“dbauth” []
  “udp”: { “user”: “foo” “password”: “bar”
}
```

} See dipper/config.py for more information

Output of fetcher: udp_variants.tsv ‘Patient’, ‘Family’, ‘Chr’, ‘Build’, ‘Chromosome Position’, ‘Reference Allele’, ‘Variant Allele’, ‘Parent of origin’, ‘Allele Type’, ‘Mutation Type’, ‘Gene’, ‘Transcript’, ‘Original Amino Acid’, ‘Variant Amino Acid’, ‘Amino Acid Change’, ‘Segregates with’, ‘Position’, ‘Exon’, ‘Inheritance model’, ‘Zygosity’, ‘dbSNP ID’, ‘1K Frequency’, ‘Number of Alleles’

udp_phenotypes.tsv ‘Patient’, ‘HPID’, ‘Present’

The script also utilizes two mapping files udp_gene_map.tsv - generated from scripts/fetch-gene-ids.py, gene symbols from udp_variants

udp_chr_rs.tsv - rsid(s) per coordinate greped from hg19 dbsnp file, then disambiguated with eutils, see scripts/dbsnp/dbsnp.py

```
UDP_SERVER = 'https://udplims-collab.nhgri.nih.gov/api'
fetch(is_dl_forced=True)
    Fetches data from udp collaboration server, see top level comments for class for more information :return:
files = {'patient_phenotypes': {'file': 'udp_phenotypes.tsv'}, 'patient_variants':
map_files = {'dbsnp_map': '../..../resources/udp/udp_chr_rs.tsv', 'gene_coord_map': '..'.
parse(limit=None)
    Override Source.parse() Args:
        :param limit (int, optional) limit the number of rows processed

Returns: :return None
```

dipper.sources.WormBase module

```
class dipper.sources.WormBase.WormBase(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source
```

This is the parser for the [C. elegans Model Organism Database (WormBase)](<http://www.wormbase.org>), from which we process genotype and phenotype data for laboratory worms (C.elegans and other nematodes).

We generate the wormbase graph to include the following information: * genes * sequence alterations (includes SNPs/del/ins/indel and

large chromosomal rearrangements)

- RNAi as expression-affecting reagents
- genotypes, and their components
- strains
- publications (and their mapping to PMIDs, if available)
- allele-to-phenotype associations (including variants by RNAi)

- genetic positional information for genes and sequence alterations

Genotypes leverage the GENO genotype model and includes both intrinsic and extrinsic genotypes. Where necessary, we create anonymous nodes of the genotype partonomy (i.e. for variant single locus complements, genomic variation complements, variant loci, extrinsic genotypes, and extrinsic genotype parts).

TODO: get people and gene expression

fetch (is_dl_forced=False)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
None

files = {'allele_pheno': {'file': 'phenotype_association.wb', 'url': 'ftp://ftp.wormbase.org/pub/wormbase/phenotype_association.wb'}}

getTestSuite()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

get_feature_type_by_class_and_biotype (ftype, biotype)

make_reagent_targeted_gene_id (gene_id, reagent_id)

parse (limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

process_allele_phenotype (limit=None)

This file compactly lists variant to phenotype associations, such that in a single row, there may be >1 variant listed per phenotype and paper. This indicates that each variant is individually associated with the given phenotype, as listed in 1+ papers. (Not that the combination of variants is producing the phenotype.) :param limit: :return:

process_disease_association (limit)

process_feature_loc (limit)

process_gene_desc (limit)

process_gene_ids (limit)

process_gene_interaction (limit)

The gene interaction file includes identified interactions, that are between two or more gene (products). In the case of interactions with >2 genes, this requires creating groups of genes that are involved in the interaction. From the wormbase help list: In the example WBInteraction000007779 it would likely be misleading to suggest that lin-12 interacts with (suppresses in this case) smo-1 ALONE or that lin-12 suppresses let-60 ALONE; the observation in the paper; see Table V in paper PMID:15990876 was that a lin-12 allele (heterozygous lin-12(n941/+)) could suppress the “multivulva” phenotype induced synthetically by simultaneous perturbation of BOTH smo-1 (by RNAi) AND let-60 (by the n2021 allele). So this is necessarily a three-gene interaction.

Therefore, we can create groups of genes based on their “status” of Effector | Effected.

Status: IN PROGRESS

Parameters limit -

Returns

process_pub_xrefs (limit=None)

process_rnai_phenotypes (limit=None)

species = '/species/c_elegans/PRJNA13758'

test_ids = {'allele': ['WBVar00087800', 'WBVar00087742', 'WBVar00144481', 'WBVar002481']}

```
update_wsnum_in_files(vernum)
    With the given version number `vernum` , update the source's version number, and replace in the file
    hashmap. the version number is in the CHECKSUMS file. :param vernum: :return:
wbdev = 'ftp://ftp.wormbase.org/pub/wormbase/releases/current-development-release'
wbprod = 'ftp://ftp.wormbase.org/pub/wormbase/releases/current-production-release'
wbrel = 'ftp://ftp.wormbase.org/pub/wormbase/releases'
```

dipper.sources.ZFIN module

```
class dipper.sources.ZFIN(graph_type, are_bnodes_skolemized)
Bases: dipper.sources.Source
```

This is the parser for the [Zebrafish Model Organism Database (ZFIN)](<http://www.zfin.org>), from which we process genotype and phenotype data for laboratory zebrafish.

We generate the zfin graph to include the following information:

- * genes
- * sequence alterations (includes SNPs/del/ins/indel and large chromosomal rearrangements)
- * transgenic constructs
- * morpholinos, talens, crisprs as expression-affecting reagents
- * genotypes, and their components
- * fish (as comprised of intrinsic and extrinsic genotypes)
- * publications (and their mapping to PMIDs, if available)
- * genotype-to-phenotype associations (including environments and stages at which they are assayed)
- * environmental components
- * orthology to human genes
- * genetic positional information for genes and sequence alterations
- * fish-to-disease model associations

Genotypes leverage the GENO genotype model and include both intrinsic and extrinsic genotypes. Where necessary, we create anonymous nodes of the genotype partonomy (such as for variant single locus complements, genomic variation complements, variant loci, extrinsic genotypes, and extrinsic genotype parts).

Furthermore, we process the genotype components to build labels in a monarch-style. This leads to genotype labels that include:

- * all genes targeted by reagents (morphants, crisprs, etc), in addition to the ones that the reagent was designed against.
- * all affected genes within deficiencies
- * complex hets being listed as gene<mutation1>/gene<mutation2> rather than gene<mutation1>/+; gene<mutation2>/+

fetch(is_dl_forced=False)

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
None

files = {'backgrounds': { 'file': 'genotype_backgrounds.txt', 'url': 'http://zfin.org'}}

getTestSuite()

An abstract method that should be overwritten with tests appropriate for the specific source. :return:

get_orthology_evidence_code(abbrev)

get_orthology_sources_from_zebrafishmine()

Fetch the zfin gene to other species orthology annotations, together with the evidence for the assertion.
Write the file locally to be read in a separate function. :return:

static make_targeted_gene_id(geneid, reagentid)

parse(limit=None)

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

process_fish(limit=None)

Fish give identifiers to the “effective genotypes” that we create. We can match these by: Fish = (intrinsic) genotype + set of morpholinos

We assume here that the intrinsic genotypes and their parts will be processed separately, prior to calling this function.

Parameters `limit` –

Returns

`process_fish_disease_models(limit=None)`

`process_orthology_evidence(limit)`

`scrub()`

Perform various data-scrubbing on the raw data files prior to parsing. For this resource, this currently includes: * remove oddities where there are “” instead of empty strings :return: None

`test_ids = {'allele': ['ZDB-ALT-010426-4', 'ZDB-ALT-010427-8', 'ZDB-ALT-011017-8', 'Z`

dipper.sources.ZFINSlim module

`class dipper.sources.ZFINSlim.ZFINSlim(graph_type, are_bnodes_skolemized)`

Bases: `dipper.sources.Source`

zfin mgf model only containing Gene to phenotype associations Using the file here: https://zfin.org/downloads/phenoGeneCleanData_fish.txt

`fetch(is_dl_forced=False)`

abstract method to fetch all data from an external resource. this should be overridden by subclasses :return: None

`files = {'g2p_clean': {'file': 'phenoGeneCleanData_fish.txt.txt', 'url': 'https://z`

`parse(limit=None)`

abstract method to parse all data from an external resource, that was fetched in fetch() this should be overridden by subclasses :return: None

dipper.utils package

Submodules

dipper.utils.CurieUtil module

`class dipper.utils.CurieUtil.CurieUtil(curie_map)`

Bases: `object`

Create compact URI

`get_base()`

`get_curie(uri)`

Get a CURIE from a URI

`get_curie_prefix(uri)`

Return the CURIE's prefix:

`get_uri(curie)`

Get a URI from a CURIE

`prefix_exists(pfx)`

dipper.utils.DipperUtil module

```
class dipper.utils.DipperUtil.DipperUtil
    Bases: object
```

Various utilities and quick methods used in this application

(A little too quick) Per: <https://www.ncbi.nlm.nih.gov/books/NBK25497/> NCBI recommends that users post

no more than three URL requests per second and limit large jobs to either weekends or between 9:00 PM and 5:00 AM Eastern time during weekdays

restructuring to make bulk queries is less likely to result in another ban for peppering them with one offs

```
static get_homologene_by_gene_num(gene_num)
```

```
static get_ncbi_id_from_symbol(gene_symbol)
```

Get ncbi gene id from symbol using monarch and mygene services :param gene_symbol: :return:

```
static get_ncbi_taxon_num_by_label(label)
```

Here we want to look up the NCBI Taxon id using some kind of label. It will only return a result if there is a unique hit.

Returns

```
static is_omim_disease(gene_id)
```

Process omim equivalencies by examining the monarch ontology scigraph As an alternative we could examine mondo.owl, since the ontology scigraph imports the output of this script which creates an odd circular dependency (even though we're querying mondo.owl through scigraph)

Parameters

- **graph** – rdfLib graph object
- **gene_id** – ncbi gene id as curie
- **omim_id** – omim id as curie

Returns

None

```
remove_control_characters(s)
```

Filters out characters in any of these unicode categories [Cc] Other, Control (65 characters)

,... [Cf] Other, Format (151 characters) [Cn] Other, Not Assigned (0 characters – none have this property) [Co] Other, Private Use (6 characters) [Cs] Other, Surrogate (6 characters)

dipper.utils.GraphUtils module

```
class dipper.utils.GraphUtils.GraphUtils(curie_map)
```

Bases: object

```
static add_property_axioms(graph, properties)
```

```
static add_property_to_graph(results, graph, property_type, property_list)
```

```
digest_id()
```

Form a deterministic digest of input Leading ‘b’ is an experiment forcing the first char to be non numeric but valid hex Not required for RDF but some other contexts do not want the leading char to be a digit

: param str wordage arbitrary string : return str

```
static get_properties_from_graph(graph)
    Wrapper for RDFLib.graph.predicates() that returns a unique set :param graph: RDFLib.graph :return: set, set of properties

write(graph, fileformat=None, file=None)
    A basic graph writer (to stdout) for any of the sources. this will write raw triples in rdFXML, unless specified. to write turtle, specify format='turtle' an optional file can be supplied instead of stdout :return: None
```

dipper.utils.TestTools module

```
class dipper.utils.TestTools
    Bases: object

    test_graph_equality(turtlish, graph)

    Parameters
        • turtlish – String of triples in turtle format without prefix header
        • graph – Graph object to test against

    Returns Boolean, True if graphs contain same set of triples
```

dipper.utils.pysed module

```
dipper.utils.pysed.replace(oldstr, newstr, infile, dryrun=False)
    Sed-like Replace function.. Usage: pysed.replace(<Old string>, <Replacement String>, <Text File>) Example: pysed.replace('xyz', 'XYZ', '/path/to/file.txt')

    This will dump the output to STDOUT instead of changing the input file. Example 'DRYRUN': pysed.replace('xyz', 'XYZ', '/path/to/file.txt', dryrun=True)

dipper.utils.pysed.rmlinematch(oldstr, infile, dryrun=False)
    Sed-like line deletion function based on given string.. Usage: pysed.rmlinematch(<Unwanted string>, <Text File>) Example: pysed.rmlinematch('xyz', '/path/to/file.txt') Example: 'DRYRUN': pysed.rmlinematch('xyz', '/path/to/file.txt', dryrun=True) This will dump the output to STDOUT instead of changing the input file.

dipper.utils.pysed.rmlinenumber(linenumber, infile, dryrun=False)
    Sed-like line deletion function based on given line number.. Usage: pysed.rmlinenumber(<Unwanted Line Number>, <Text File>) Example: pysed.rmlinenumber(10, '/path/to/file.txt') Example 'DRYRUN': pysed.rmlinenumber(10, '/path/to/file.txt', dryrun=True) #This will dump the output to STDOUT instead of changing the input file.
```

dipper.utils.romanplus module

```
exception dipper.utils.romanplus.InvalidRomanNumeralError
    Bases: dipper.utils.romanplus.RomanError

exception dipper.utils.romanplus.NotIntegerError
    Bases: dipper.utils.romanplus.RomanError

exception dipper.utils.romanplus.OutOfRangeError
    Bases: dipper.utils.romanplus.RomanError

exception dipper.utils.romanplus.RomanError
    Bases: Exception
```

```
dipper.utils.romanplus.fromRoman(s)
```

convert Roman numeral to integer

```
dipper.utils.romanplus.toRoman(n)
```

convert integer to Roman numeral

Submodules

dipper.config module

```
dipper.config.conf = {'keys': {'omim': ''}}
```

Load the configuration file ‘conf.json’, if it exists. It isn’t always required, but may be for some sources.

conf.json may contain sensitive info and should not live in a public repo

```
dipper.config.get_config()
```

dipper.curie_map module

Acroname central

Load the curie mapping file ‘curie_map.yaml’, it is necessary for most resources

```
dipper.curie_map.get()
```

```
dipper.curie_map.get_base()
```

3.2 Source APIs

CHAPTER 4

Indices and tables

- genindex
- modindex
- search

Python Module Index

d

dipper, 13
dipper.config, 63
dipper.curie_map, 63
dipper.graph, 13
dipper.graph.Graph, 13
dipper.graph.RDFGraph, 13
dipper.graph.StreamedGraph, 14
dipper.models, 14
dipper.models.assoc, 14
dipper.models.assoc.Association, 14
dipper.models.assoc.Chem2DiseaseAssoc, 15
dipper.models.assoc.D2PAssoc, 16
dipper.models.assocDispositionAssoc, 16
dipper.models.assoc.G2PAssoc, 16
dipper.models.assoc.InteractionAssoc, 17
dipper.models.assoc.OrthologyAssoc, 17
dipper.models.Dataset, 18
dipper.models.Environment, 18
dipper.models.Evidence, 19
dipper.models.Family, 20
dipper.models.GenomicFeature, 20
dipper.models.Genotype, 22
dipper.models.Model, 25
dipper.models.Pathway, 26
dipper.models.Provenance, 27
dipper.models.Reference, 27
dipper.sources, 28
dipper.sources.AnimalQTlDb, 28
dipper.sources.Bgee, 29
dipper.sources.BioGrid, 30
dipper.sources.ClinVar, 31
dipper.sources.Coriell, 31
dipper.sources.CTD, 30
dipper.sources.Decipher, 32
dipper.sources.Ensembl, 33
dipper.sources.EOM, 33
dipper.sources.FlyBase, 34
dipper.sources.GeneOntology, 35
dipper.sources.GeneReviews, 36
dipper.sources.GWASCatalog, 34
dipper.sources.HGNC, 36
dipper.sources.HPOAnnotations, 37
dipper.sources.IMPC, 38
dipper.sources.KEGG, 39
dipper.sources.MGI, 40
dipper.sources.MGISlim, 40
dipper.sources.MMRRC, 41
dipper.sources.Monarch, 42
dipper.sources.Monochrom, 43
dipper.sources.MPD, 41
dipper.sources.MyChem, 44
dipper.sources.MyDrug, 44
dipper.sources.NCBIGene, 45
dipper.sources.OMA, 46
dipper.sources.OMIA, 46
dipper.sources.OMIM, 48
dipper.sources.Orphanet, 49
dipper.sources.OrthoXML, 49
dipper.sources.Panther, 50
dipper.sources.PostgreSQLSource, 50
dipper.sources.Reactome, 51
dipper.sources.RGD, 51
dipper.sources.SGD, 52
dipper.sources.Source, 52
dipper.sources.StringDB, 55
dipper.sources.UCSCBands, 55
dipper.sources.UDP, 56
dipper.sources.WormBase, 57
dipper.sources.ZFIN, 59
dipper.sources.ZFINSlim, 60
dipper.utils, 60
dipper.utils.CurieUtil, 60
dipper.utils.DipperUtil, 61
dipper.utils.GraphUtils, 61
dipper.utils.pysed, 62

`dipper.utils.romanplus`, 62
`dipper.utils.TestTools`, 62

Index

A

- add_agent_to_graph() (dip-
per.models.Provenance.Provenance method), 27
- add_assay_to_graph() (dip-
per.models.Provenance.Provenance method), 27
- add_assertion() (dipper.models.Provenance.Provenance method), 27
- add_association_to_graph() (dip-
per.models.assoc.Association.Assoc method), 14
- add_association_to_graph() (dip-
per.models.assoc.D2PAssoc.D2PAssoc method), 16
- add_association_to_graph() (dip-
per.models.assoc.G2PAssoc.G2PAssoc method), 17
- add_common_files_to_file_list() (dip-
per.sources.HPOAnnotations.HPOAnnotations method), 37
- add_data_individual() (dipper.models.Evidence.Evidence method), 19
- add_date() (dipper.models.assoc.Association.Assoc method), 14
- add_date_created() (dip-
per.models.Provenance.Provenance method), 27
- add_evidence() (dipper.models.assoc.Association.Assoc method), 14
- add_evidence() (dipper.models.Evidence.Evidence method), 19
- add_gene_family_to_graph() (dip-
per.models.assoc.OntologyAssoc.OntologyAssoc method), 17
- add_orthologs_by_gene_group() (dip-
per.sources.NCBIGene.NCBIGene method), 45
- add_predicate_object() (dip-
- per.models.assoc.Association.Assoc method), 14
- add_property_axioms() (dip-
per.utils.GraphUtils.GraphUtils static method), 61
- add_property_to_graph() (dip-
per.utils.GraphUtils.GraphUtils static method), 61
- add_protein_to_graph (dip-
per.sources.OrthoXML.OrthoXML attribute), 49
- add_provenance() (dip-
per.models.assoc.Association.Assoc method), 14
- add_relation() (dipper.sources.MyChem.MyChem static method), 44
- add_source() (dipper.models.assoc.Association.Assoc method), 14
- add_source() (dipper.models.Evidence.Evidence method), 19
- add_study_measure() (dip-
per.models.Provenance.Provenance method), 27
- add_study_parts() (dip-
per.models.Provenance.Provenance method), 27
- add_study_to_measurements() (dip-
per.models.Provenance.Provenance method), 27
- add_supporting_data() (dip-
per.models.Evidence.Evidence method), 20
- add_supporting_evidence() (dip-
per.models.Evidence.Evidence method), 20
- add_supporting_publication() (dip-
per.models.Evidence.Evidence method), 20
- addAffectedLocus() (dipper.models.Genotype.Genotype method), 22

addAllele() (dipper.models.Genotype.Genotype method), 22	addFeatureStartLocation() (dipper.models.GenomicFeature.Feature method), 21
addAlleleOfGene() (dipper.models.Genotype.Genotype method), 22	addFeatureToGraph() (dipper.models.GenomicFeature.Feature method), 21
addAuthor() (dipper.models.Reference.Reference method), 28	addGene() (dipper.models.Genotype.Genotype method), 23
addBlankNodeAnnotation() (dipper.models.Model.Model method), 25	addGeneProduct() (dipper.models.Genotype.Genotype method), 23
addChromosome() (dipper.models.Genotype.Genotype method), 22	addGeneTargetingReagent() (dipper.models.Genotype.Genotype method), 23
addChromosomeClass() (dipper.models.Genotype.Genotype method), 22	addGeneTargetingReagentToGenotype() (dipper.models.Genotype.Genotype method), 23
addChromosomeInstance() (dipper.models.Genotype.Genotype method), 23	addGeneToPathway() (dipper.models.Pathway.Pathway method), 26
addClassToGraph() (dipper.models.Model.Model method), 25	addGenome() (dipper.models.Genotype.Genotype method), 23
addComment() (dipper.models.Model.Model method), 25	addGenomicBackground() (dipper.models.Genotype.Genotype method), 23
addComponentAttributes() (dipper.models.Environment.Environment method), 19	addGenomicBackgroundToGenotype() (dipper.models.Genotype.Genotype method), 23
addComponentToEnvironment() (dipper.models.Environment.Environment method), 19	addGenotype() (dipper.models.Genotype.Genotype method), 23
addComponentToPathway() (dipper.models.Pathway.Pathway method), 26	addIndividualToGraph() (dipper.models.Model.Model method), 26
addConstruct() (dipper.models.Genotype.Genotype method), 23	addLabel() (dipper.models.Model.Model method), 26
addDefinition() (dipper.models.Model.Model method), 25	addMember() (dipper.models.Family.Family method), 20
addDepiction() (dipper.models.Model.Model method), 25	addMemberOf() (dipper.models.Family.Family method), 20
addDeprecatedClass() (dipper.models.Model.Model method), 25	addMemberOfPopulation() (dipper.models.Genotype.Genotype method), 23
addDeprecatedIndividual() (dipper.models.Model.Model method), 25	addOntologyDeclaration() (dipper.models.Model.Model method), 26
addDerivesFrom() (dipper.models.Genotype.Genotype method), 23	addOWLPropertyClassRestriction() (dipper.models.Model.Model method), 26
addDescription() (dipper.models.Model.Model method), 25	addOWLVersionInfo() (dipper.models.Model.Model method), 26
addEnvironment() (dipper.models.Environment.Environment method), 19	addOWLVersionIRI() (dipper.models.Model.Model method), 26
addEnvironmentalCondition() (dipper.models.Environment.Environment method), 19	addPage() (dipper.models.Reference.Reference method), 28
addEquivalentClass() (dipper.models.Model.Model method), 25	addParts() (dipper.models.Genotype.Genotype method), 23
addFeatureEndLocation() (dipper.models.GenomicFeature.Feature method), 21	addPartsToVSLC() (dipper.models.Genotype.Genotype method), 23
addFeatureProperty() (dipper.models.GenomicFeature.Feature method), 21	addPathway() (dipper.models.Pathway.Pathway method), 27
	addPerson() (dipper.models.Model.Model method), 26

addPolypeptide() (dipper.models.Genotype.Genotype method), 23
 addPositionToGraph() (dipper.models.GenomicFeature.Feature method), 21
 addReagentTargetedGene() (dipper.models.Genotype.Genotype method), 24
 addReferenceGenome() (dipper.models.Genotype.Genotype method), 24
 addRefToGraph() (dipper.models.Reference.Reference method), 28
 addRegionPositionToGraph() (dipper.models.GenomicFeature.Feature method), 21
 addSameIndividual() (dipper.models.Model.Model method), 26
 addSequenceAlteration() (dipper.models.Genotype.Genotype method), 24
 addSequenceAlterationToVariantLocus() (dipper.models.Genotype.Genotype method), 24
 addSequenceDerivesFrom() (dipper.models.Genotype.Genotype method), 24
 addSubClass() (dipper.models.Model.Model method), 26
 addSubsequenceOfFeature() (dipper.models.GenomicFeature.Feature method), 21
 addSynonym() (dipper.models.Model.Model method), 26
 addTargetedGeneComplement() (dipper.models.Genotype.Genotype method), 24
 addTargetedGeneSubregion() (dipper.models.Genotype.Genotype method), 24
 addTaxon() (dipper.models.Genotype.Genotype method), 24
 addTaxonToFeature() (dipper.models.GenomicFeature.Feature method), 21
 addTitle() (dipper.models.Reference.Reference method), 28
 addTriple() (dipper.graph.Graph.Graph method), 13
 addTriple() (dipper.graph.RDFGraph.RDFGraph method), 13
 addTriple() (dipper.graph.StreamedGraph.StreamedGraph method), 14
 addTriple() (dipper.models.Model.Model method), 26
 addType() (dipper.models.Model.Model method), 26
 addVSLCtoParent() (dipper.models.Genotype.Genotype method), 24
 addXref() (dipper.models.Model.Model method), 26
 AnimalQTLdb (class in dipper.sources.AnimalQTLdb), 28
 annotation_properties (dipper.models.assoc.Association.Assoc attribute), 14
 annotation_properties (dipper.models.Environment.Environment attribute), 19
 annotation_properties (dipper.models.GenomicFeature.Feature attribute), 22
 annotation_properties (dipper.models.Genotype.Genotype attribute), 24
 annotation_properties (dipper.models.Model.Model attribute), 26
 annotation_properties (dipper.models.Reference.Reference attribute), 28
 Assoc (class in dipper.models.assoc.Association), 14
 assoc_types (dipper.models.assoc.Association.Assoc attribute), 15

B

BENCHMARK_BASE (dipper.sources.OMA.OMA attribute), 46
 Bgee (class in dipper.sources.Bgee), 29
 BGEE_FTP (dipper.sources.Bgee.Bgee attribute), 29
 bind_all_namespaces() (dipper.graph.RDFGraph.RDFGraph method), 13
 BioGrid (class in dipper.sources.BioGrid), 30
 biogrid_ids (dipper.sources.BioGrid.BioGrid attribute), 30
 build_measurement_description() (dipper.sources.MPD.MPD static method), 42

C

check_header() (dipper.sources.MPD.MPD static method), 42
 check_uniprot() (dipper.sources.MyChem.MyChem static method), 44
 checkIfRemoteIsNewer() (dipper.sources.Bgee.Bgee method), 29
 checkIfRemoteIsNewer() (dipper.sources.MyDrug.MyDrug method), 44
 checkIfRemoteIsNewer() (dipper.sources.Source.Source method), 52
 Chem2DiseaseAssoc (class in dipper.models.assoc.Chem2DiseaseAssoc), 15
 chunks() (dipper.sources.MyChem.MyChem static method), 44

clean_db_prefix() (dipper.sources.GeneOntology.GeneOntology method), 35
clean_protein_id() (dipper.sources.OrthoXML.OrthoXML method), 49
clean_up_omim_genes() (dipper.sources.OMIA.OMIA method), 47
ClinVar (class in dipper.sources.ClinVar), 31
compare_checksums() (dipper.sources.IMPC.IMPC method), 39
compare_local_remote_bytes() (dipper.sources.Source.Source method), 53
conf (in module dipper.config), 63
Coriell (class in dipper.sources.Coriell), 31
create_books() (dipper.sources.GeneReviews.GeneReviews method), 36
CTD (class in dipper.sources.CTD), 30
curie_util (dipper.graph.RDFGraph.RDFGraph attribute), 13
curie_util (dipper.graph.StreamedGraph.StreamedGraph attribute), 14
CurieUtil (class in dipper.utils.CurieUtil), 60

D

d2p_object_properties (dipper.models.assoc.D2PAssoc.D2PAssoc attribute), 16
D2PAssoc (class in dipper.models.assoc.D2PAssoc), 16
data_properties (dipper.models.GenomicFeature.Feature attribute), 22
data_property (dipper.models.Evidence.Evidence attribute), 20
data_types (dipper.models.Evidence.Evidence attribute), 20
Dataset (class in dipper.models.Dataset), 18
datatype_properties (dipper.models.assoc.Association.Assoc attribute), 15
datatype_properties (dipper.models.Model.Model attribute), 26
Decipher (class in dipper.sources.Decipher), 32
declareAsOntology() (dipper.sources.Source.Source method), 53
default_node_list() (dipper.sources.OrthoXML.OrthoXMLParser method), 50
DEFAULT_TAXA (dipper.sources.Bgee.Bgee attribute), 29
DEFAULT_TAXA (dipper.sources.StringDB.StringDB attribute), 55
digest_id() (dipper.utils.GraphUtils.GraphUtils method), 61
dipper (module), 13
dipper.config (module), 63
dipper.curie_map (module), 63
dipper.graph (module), 13
dipper.graph.Graph (module), 13
dipper.graph.RDFGraph (module), 13
dipper.graph.StreamedGraph (module), 14
dipper.models (module), 14
dipper.models.assoc (module), 14
dipper.models.assoc.Association (module), 14
dipper.models.assoc.Chem2DiseaseAssoc (module), 15
dipper.models.assoc.D2PAssoc (module), 16
dipper.models.assocDispositionAssoc (module), 16
dipper.models.assoc.G2PAssoc (module), 16
dipper.models.assoc.InteractionAssoc (module), 17
dipper.models.assoc.OrthologyAssoc (module), 17
dipper.models.Dataset (module), 18
dipper.models.Environment (module), 18
dipper.models.Evidence (module), 19
dipper.models.Family (module), 20
dipper.models.GenomicFeature (module), 20
dipper.models.Genotype (module), 22
dipper.models.Model (module), 25
dipper.models.Pathway (module), 26
dipper.models.Provenance (module), 27
dipper.models.Reference (module), 27
dipper.sources (module), 28
dipper.sources.AnimalQTLdb (module), 28
dipper.sources.Bgee (module), 29
dipper.sources.BioGrid (module), 30
dipper.sources.ClinVar (module), 31
dipper.sources.Coriell (module), 31
dipper.sources.CTD (module), 30
dipper.sources.Decipher (module), 32
dipper.sources.Ensembl (module), 33
dipper.sources.EOM (module), 33
dipper.sources.FlyBase (module), 34
dipper.sources.GeneOntology (module), 35
dipper.sources.GeneReviews (module), 36
dipper.sources.GWASCatalog (module), 34
dipper.sources.HGNC (module), 36
dipper.sources.HPOAnnotations (module), 37
dipper.sources.IMPC (module), 38
dipper.sources.KEGG (module), 39
dipper.sources.MGI (module), 40
dipper.sources.MGISlim (module), 40
dipper.sources.MMRRC (module), 41
dipper.sources.Monarch (module), 42
dipper.sources.Monochrom (module), 43
dipper.sources.MPD (module), 41
dipper.sources.MyChem (module), 44
dipper.sources.MyDrug (module), 44
dipper.sources.NCBIGene (module), 45
dipper.sources.OMA (module), 46
dipper.sources.OMIA (module), 46

dipper.sources.OMIM (module), 48
dipper.sources.Orphanet (module), 49
dipper.sources.OrthoXML (module), 49
dipper.sources.Panther (module), 50
dipper.sources.PostgreSQLSource (module), 50
dipper.sources.Reactome (module), 51
dipper.sources.RGD (module), 51
dipper.sources.SGD (module), 52
dipper.sources.Source (module), 52
dipper.sources.StringDB (module), 55
dipper.sources.UCSCBands (module), 55
dipper.sources.UDP (module), 56
dipper.sources.WormBase (module), 57
dipper.sources.ZFIN (module), 59
dipper.sources.ZFINSlim (module), 60
dipper.utils (module), 60
dipper.utils.CurieUtil (module), 60
dipper.utils.DipperUtil (module), 61
dipper.utils.GraphUtils (module), 61
dipper.utils.pysed (module), 62
dipper.utils.romanplus (module), 62
dipper.utils.TestTools (module), 62
DipperUtil (class in dipper.utils.DipperUtil), 61
DispositionAssoc (class in dipper.models.assocDispositionAssoc), 16

E

eco_dict (dipper.sources.HPOAnnotations.HPOAnnotations attribute), 37
Ensembl (class in dipper.sources.Ensembl), 33
Environment (class in dipper.models.Environment), 18
environment_parts (dipper.models.Environment.Environment attribute), 19
EOM (class in dipper.sources.EOM), 33
Evidence (class in dipper.models.Evidence), 19
evidence_types (dipper.models.Evidence.Evidence attribute), 20
execute_query() (dipper.sources.MyChem.MyChem static method), 44
extract_pairwise_relations() (dipper.sources.OrthoXML.OrthoXMLParser method), 50
extract_taxon_info() (dipper.sources.OrthoXML.OrthoXML method), 49

F

Family (class in dipper.models.Family), 20
Feature (class in dipper.models.GenomicFeature), 20
fetch() (dipper.sources.AnimalQTLdb.AnimalQTLdb method), 29
fetch() (dipper.sources.Bgee.Bgee method), 29
fetch() (dipper.sources.BioGrid.BioGrid method), 30

fetch() (dipper.sources.ClinVar.ClinVar method), 31
fetch() (dipper.sources.Coriell.Coriell method), 32
fetch() (dipper.sources.CTD.CTD method), 30
fetch() (dipper.sources.Decipher.Decipher method), 32
fetch() (dipper.sources.Ensembl.Ensembl method), 33
fetch() (dipper.sources.EOM.EOM method), 33
fetch() (dipper.sources.FlyBase.FlyBase method), 34
fetch() (dipper.sources.GeneOntology.GeneOntology method), 35
fetch() (dipper.sources.GeneReviews.GeneReviews method), 36
fetch() (dipper.sources.GWASCatalog.GWASCatalog method), 35
fetch() (dipper.sources.HGNC.HGNC method), 37
fetch() (dipper.sources.HPOAnnotations.HPOAnnotations method), 37
fetch() (dipper.sources.IMPC.IMPC method), 39
fetch() (dipper.sources.KEGG.KEGG method), 39
fetch() (dipper.sources.MGI.MGI method), 40
fetch() (dipper.sources.MGISlim.MGISlim method), 41
fetch() (dipper.sources.MMRRC.MMRRC method), 41
fetch() (dipper.sources.Monarch.Monarch method), 42
fetch() (dipper.sources.Monochrom.Monochrom method), 43
fetch() (dipper.sources.MPD.MPD method), 42
fetch() (dipper.sources.MyChem.MyChem method), 44
fetch() (dipper.sources.MyDrug.MyDrug method), 45
fetch() (dipper.sources.NCBIGene.NCBIGene method), 46
fetch() (dipper.sources.OMIA.OMIA method), 47
fetch() (dipper.sources.OMIM.OMIM method), 48
fetch() (dipper.sources.Orphanet.Orphanet method), 49
fetch() (dipper.sources.OrthoXML.OrthoXML method), 49
fetch() (dipper.sources.Panther.Panther method), 50
fetch() (dipper.sources.Reactome.Reactome method), 51
fetch() (dipper.sources.RGD.RGD method), 51
fetch() (dipper.sources.SGD.SGD method), 52
fetch() (dipper.sources.Source.Source method), 53
fetch() (dipper.sources.StringDB.StringDB method), 55
fetch() (dipper.sources.UCSCBands.UCSCBands method), 56
fetch() (dipper.sources.UDP.UDP method), 57
fetch() (dipper.sources.WormBase.WormBase method), 58
fetch() (dipper.sources.ZFIN.ZFIN method), 59
fetch() (dipper.sources.ZFINSlim.ZFINSlim method), 60
fetch_from_mychem() (dipper.sources.MyChem.MyChem method), 44
fetch_from_pgdb() (dipper.sources.PostgreSQLSource.PostgreSQLSource method), 50

fetch_from_url() (dipper.sources.Source.Source method), 53
fetch_protein_gene_map() (dipper.sources.Ensembl.Ensembl method), 33
fetch_protein_list() (dipper.sources.Ensembl.Ensembl method), 33
fetch_query_from_pgdb() (dipper.sources.PostgreSQLSource.PostgreSQLSource method), 51
fetch_transgene_genes_from_db() (dipper.sources.MGI.MGI method), 40
fetch_uniprot_gene_map() (dipper.sources.Ensembl.Ensembl method), 33
file_len() (dipper.sources.Source.Source method), 53
files (dipper.sources.AnimalQTLdb.AnimalQTLdb attribute), 29
files (dipper.sources.Bgee.Bgee attribute), 29
files (dipper.sources.BioGrid.BioGrid attribute), 30
files (dipper.sources.ClinVar.ClinVar attribute), 31
files (dipper.sources.Coriell.Coriell attribute), 32
files (dipper.sources.CTD.CTD attribute), 30
files (dipper.sources.Decipher.Decipher attribute), 32
files (dipper.sources.Ensembl.Ensembl attribute), 33
files (dipper.sources.EOM.EOM attribute), 33
files (dipper.sources.FlyBase.FlyBase attribute), 34
files (dipper.sources.GeneOntology.GeneOntology attribute), 35
files (dipper.sources.GeneReviews.GeneReviews attribute), 36
files (dipper.sources.GWASCatalog.GWASCatalog attribute), 35
files (dipper.sources.HGNC.HGNC attribute), 37
files (dipper.sources.HPOAnnotations.HPOAnnotations attribute), 37
files (dipper.sources.IMPC.IMPC attribute), 39
files (dipper.sources.KEGG.KEGG attribute), 39
files (dipper.sources.MMRRC.MMRRC attribute), 41
files (dipper.sources.Monochrom.Monochrom attribute), 43
files (dipper.sources.MPD.MPD attribute), 42
files (dipper.sources.MyDrug.MyDrug attribute), 45
files (dipper.sources.NCBIGene.NCBIGene attribute), 46
files (dipper.sources.OMA.OMA attribute), 46
files (dipper.sources.OMIA.OMIA attribute), 47
files (dipper.sources.OMIM.OMIM attribute), 48
files (dipper.sources.Orphanet.Orphanet attribute), 49
files (dipper.sources.OrthoXML.OrthoXML attribute), 49
files (dipper.sources.Panther.Panther attribute), 50
files (dipper.sources.Reactome.Reactome attribute), 52
files (dipper.sources.RGD.RGD attribute), 51
files (dipper.sources.SGD.SGD attribute), 52
files (dipper.sources.Source.Source attribute), 53
files (dipper.sources.UCSCBands.UCSCBands attribute), 56
files (dipper.sources.UDP.UDP attribute), 57
files (dipper.sources.WormBase.WormBase attribute), 58
files (dipper.sources.ZFIN.ZFIN attribute), 59
files (dipper.sources.ZFINSlim.ZFINSlim attribute), 60
filter_keep_phenotype_entry_ids() (in module dipper.sources.OMIM), 48
FlyBase (class in dipper.sources.FlyBase), 34
format_actions() (dipper.sources.MyChem.MyChem static method), 44
fromRoman() (in module dipper.utils.romanplus), 62

G

g2p_types (dipper.models.assoc.G2PAssoc.G2PAssoc attribute), 17
G2PAssoc (class in dipper.models.assoc.G2PAssoc), 16
GeneOntology (class in dipper.sources.GeneOntology), 35
GeneReviews (class in dipper.sources.GeneReviews), 36
genoparts (dipper.models.Genotype.Genotype attribute), 24
Genotype (class in dipper.models.Genotype), 22
get() (in module dipper.curie_map), 63
get_association_id() (dipper.models.assoc.Association.Assoc method), 15
get_base() (dipper.utils.CurieUtil.CurieUtil method), 60
get_base() (in module dipper.curie_map), 63
get_children() (dipper.sources.OrthoXML.OrthoXMLParser method), 50
get_common_files() (dipper.sources.HPOAnnotations.HPOAnnotations method), 37
get_config() (in module dipper.config), 63
get_curie() (dipper.utils.CurieUtil.CurieUtil method), 60
get_curie_prefix() (dipper.utils.CurieUtil.CurieUtil method), 60
get_doid_ids_for_unpadding() (dipper.sources.HPOAnnotations.HPOAnnotations method), 37
get_drug_record() (dipper.sources.MyChem.MyChem static method), 44
get_eco_map() (dipper.sources.Source.Source static method), 53
get_feature_type_by_class_and_biotype() (dipper.sources.WormBase.WormBase method), 58
get_file_md5() (dipper.sources.Source.Source method), 53
get_files() (dipper.sources.Source.Source method), 53
get_homologene_by_gene_num() (dipper.utils.DipperUtil.DipperUtil static method), 61
get_inchikeys() (dipper.sources.MyChem.MyChem static method), 44

get_license() (dipper.models.Dataset.Dataset method), 18
 get_local_file_size() (dipper.sources.Source.Source method), 53
 get_ncbi_id_from_symbol() (dipper.utils.DipperUtil.DipperUtil static method), 61
 get_ncbi_taxon_num_by_label() (dipper.utils.DipperUtil.DipperUtil static method), 61
 get_omim_id_from_entry() (in module dipper.sources.OMIM), 48
 get_orthology_evidence_code() (dipper.sources.ZFIN.ZFIN method), 59
 get_orthology_sources_from_zebrafishmine() (dipper.sources.ZFIN.ZFIN method), 59
 get_properties() (dipper.models.assoc.Association.Assoc method), 15
 get_properties_from_graph() (dipper.utils.GraphUtils.GraphUtils static method), 61
 get_remote_content_len() (dipper.sources.Source.Source method), 54
 get_symbol_id_map() (dipper.sources.HGNC.HGNC method), 37
 get_uniprot_entrez_id_map() (dipper.sources.GeneOntology.GeneOntology method), 35
 get_uri() (dipper.utils.CurieUtil.CurieUtil method), 60
 getChrPartTypeByNotation() (in module dipper.sources.Monochrom), 44
 getGraph() (dipper.models.Dataset.Dataset method), 18
 getTestSuite() (dipper.sources.AnimalQTLdb.AnimalQTLdb method), 29
 getTestSuite() (dipper.sources.BioGrid.BioGrid method), 30
 getTestSuite() (dipper.sources.ClinVar.ClinVar method), 31
 getTestSuite() (dipper.sources.Coriell.Coriell method), 32
 getTestSuite() (dipper.sources.CTD.CTD method), 30
 getTestSuite() (dipper.sources.Ensembl.Ensembl method), 34
 getTestSuite() (dipper.sources.EOM.EOM method), 33
 getTestSuite() (dipper.sources.FlyBase.FlyBase method), 34
 getTestSuite() (dipper.sources.GeneOntology.GeneOntology method), 35
 getTestSuite() (dipper.sources.GeneReviews.GeneReviews method), 36
 getTestSuite() (dipper.sources.GWASCatalog.GWASCatalog method), 35
 getTestSuite() (dipper.sources.HGNC.HGNC method), 37
 getTestSuite() (dipper.sources.HPOAnnotations.HPOAnnotations method), 37
 getTestSuite() (dipper.sources.IMPC.IMPC method), 39
 getTestSuite() (dipper.sources.KEGG.KEGG method), 39
 getTestSuite() (dipper.sources.MGI.MGI method), 40
 getTestSuite() (dipper.sources.MMRRC.MMRRC method), 41
 getTestSuite() (dipper.sources.Monochrom.Monochrom method), 43
 getTestSuite() (dipper.sources.MPD.MPD method), 42
 getTestSuite() (dipper.sources.NCBIGene.NCBIGene method), 46
 getTestSuite() (dipper.sources.OMA.OMA method), 46
 getTestSuite() (dipper.sources.OMIA.OMIA method), 47
 getTestSuite() (dipper.sources.OMIM.OMIM method), 48
 getTestSuite() (dipper.sources.Orphanet.Orphanet method), 49
 getTestSuite() (dipper.sources.Panther.Panther method), 50
 getTestSuite() (dipper.sources.Source.Source method), 53
 getTestSuite() (dipper.sources.UCSCBands.UCSCBands method), 56
 getTestSuite() (dipper.sources.WormBase.WormBase method), 58
 getTestSuite() (dipper.sources.ZFIN.ZFIN method), 59
 Graph (class in dipper.graph.Graph), 13
 GraphUtils (class in dipper.utils.GraphUtils), 61
 GWASCatalog (class in dipper.sources.GWASCatalog), 34
 GWASFILE (dipper.sources.GWASCatalog.GWASCatalog attribute), 35
 GWASFTP (dipper.sources.GWASCatalog.GWASCatalog attribute), 35

H

hash_id() (dipper.sources.Source.Source static method), 54
 HGGP (dipper.sources.UCSCBands.UCSCBands attribute), 56
 HGNC (class in dipper.sources.HGNC), 36
 HPOAnnotations (class in dipper.sources.HPOAnnotations), 37

I

IMPC (class in dipper.sources.IMPC), 38
 interaction_object_properties (dipper.models.assoc.InteractionAssoc.InteractionAssoc attribute), 17
 InteractionAssoc (class in dipper.models.assoc.InteractionAssoc), 17
 InvalidRomanNumeralError, 62
 is_internal_node() (dipper.sources.OrthoXML.OrthoXMLParser method), 50

is_leaf() (dipper.sources.OrthoXML.OrthoXMLParser method), 50
is_omim_disease() (dipper.utils.DipperUtil.DipperUtil static method), 61

K

KEGG (class in dipper.sources.KEGG), 39

L

leaf_label() (dipper.sources.OrthoXML.OrthoXMLParser method), 50

M

make_allele_by_consequence() (dipper.sources.Decipher.Decipher method), 32
make_apo_map() (dipper.sources.SGD.SGD static method), 52
make_association() (dipper.sources.RGD.RGD method), 51
make_association() (dipper.sources.SGD.SGD method), 52
make_association_id() (dipper.models.assoc.Association.Assoc method), 15
make_breed_id() (dipper.sources.OMIA.OMIA method), 47
make_c2p_assoc_id() (dipper.models.assoc.Chem2DiseaseAssoc.Chem2DiseaseAssoc method), 15
make_d2p_id() (dipper.models.assoc.D2PAssoc.D2PAssoc method), 16
make_experimental_model_with_genotype() (dipper.models.Genotype.Genotype method), 24
make_g2p_id() (dipper.models.assoc.G2PAssoc.G2PAssoc method), 17
make_id() (dipper.sources.Source.Source static method), 54
make_parent_bands() (dipper.sources.Monochrom.Monochrom method), 43
make_reagent_targeted_gene_id() (dipper.sources.WormBase.WormBase method), 58
make_targeted_gene_id() (dipper.sources.ZFIN.ZFIN static method), 59
make_triples() (dipper.sources.MyChem.MyChem method), 44
make_variant_locus_label() (dipper.models.Genotype.Genotype method), 24
make_vscl_label() (dipper.models.Genotype.Genotype method), 24

makeChromID() (in module dipper.models.GenomicFeature), 22
makeChromLabel() (in module dipper.models.GenomicFeature), 22
makeGenomeID() (dipper.models.Genotype.Genotype method), 24

makeLeader() (dipper.models.Model.Model method), 26
map_files (dipper.sources.GeneOntology.GeneOntology attribute), 35
map_files (dipper.sources.IMPC.IMPC attribute), 39
map_files (dipper.sources.Reactome.Reactome attribute), 52

map_files (dipper.sources.UDP.UDP attribute), 57
map_omnia_group_category_to_ontology_id() (dipper.sources.OMIA.OMIA method), 47
map_type_of_gene() (dipper.sources.NCBIGene.NCBIGene static method), 46
map_type_of_region() (dipper.sources.Monochrom.Monochrom method), 43
mgd_agent_id (dipper.sources.MPD.MPD attribute), 42
mgd_agent_label (dipper.sources.MPD.MPD attribute), 42
mgd_agent_type (dipper.sources.MPD.MPD attribute), 42

MGI (class in dipper.sources.MGI), 40
MGISlim (class in dipper.sources.MGISlim), 40
MMRRC (class in dipper.sources.MMRRC), 41
Model (class in dipper.models.Model), 25
Monarch (class in dipper.sources.Monarch), 42
Monochrom (class in dipper.sources.Monochrom), 43
MPD (class in dipper.sources.MPD), 41
MPDDL (dipper.sources.MPD.MPD attribute), 42
MY_DRUG_API (dipper.sources.MyDrug.MyDrug attribute), 44
MyChem (class in dipper.sources.MyChem), 44
MyDrug (class in dipper.sources.MyDrug), 44

N

namespaces (dipper.sources.Source.Source attribute), 54
NCBIGene (class in dipper.sources.NCBIGene), 45
normalise_units() (dipper.sources.MPD.MPD static method), 42
NotIntegerError, 62

O

object_properties (dipper.models.assoc.Association.Assoc attribute), 15
object_properties (dipper.models.Environment.Environment attribute), 19
object_properties (dipper.models.Evidence.Evidence attribute), 20

object_properties (dipper.models.Family.Family attribute), 20
 object_properties (dipper.models.GenomicFeature.Feature attribute), 22
 object_properties (dipper.models.Genotype.Genotype attribute), 24
 object_properties (dipper.models.Model.Model attribute), 26
 object_properties (dipper.models.Pathway.Pathway attribute), 27
 object_properties (dipper.models.Provenance.Provenance attribute), 27
 OMA (class in dipper.sources.OMA), 46
 OMIA (class in dipper.sources.OMIA), 46
 OMIM (class in dipper.sources.OMIM), 48
 open_and_parse_yaml() (dipper.sources.Source.Source static method), 54
 Orphanet (class in dipper.sources.Orphanet), 49
 ortho_rel (dipper.models.assoc.OntologyAssoc.Ontology attribute), 18
 OrthologyAssoc (class in dipper.models.assoc.OntologyAssoc), 17
 OrthoXML (class in dipper.sources.OrthoXML), 49
 OrthoXMLParser (class in dipper.sources.OrthoXML), 49
 OutOfRangeError, 62

P

Panther (class in dipper.sources.Panther), 50
 parse() (dipper.sources.AnimalQTLdb.AnimalQTLdb method), 29
 parse() (dipper.sources.Bgee.Bgee method), 29
 parse() (dipper.sources.BioGrid.BioGrid method), 30
 parse() (dipper.sources.ClinVar.ClinVar method), 31
 parse() (dipper.sources.Coriell.Coriell method), 32
 parse() (dipper.sources.CTD.CTD method), 31
 parse() (dipper.sources.Decipher.Decipher method), 33
 parse() (dipper.sources.Ensembl.Ensembl method), 34
 parse() (dipper.sources.EOM.EOM method), 33
 parse() (dipper.sources.FlyBase.FlyBase method), 34
 parse() (dipper.sources.GeneOntology.GeneOntology method), 35
 parse() (dipper.sources.GeneReviews.GeneReviews method), 36
 parse() (dipper.sources.GWASCatalog.GWASCatalog method), 35
 parse() (dipper.sources.HGNC.HGNC method), 37
 parse() (dipper.sources.HPOAnnotations.HPOAnnotations method), 38
 parse() (dipper.sources.IMPC.IMPC method), 39
 parse() (dipper.sources.KEGG.KEGG method), 39
 parse() (dipper.sources.MGI.MGI method), 40
 parse() (dipper.sources.MGISlim.MGISlim method), 41
 parse() (dipper.sources.MMRC.MMRC method), 41
 parse() (dipper.sources.Monarch.Monarch method), 42
 parse() (dipper.sources.Monochrom.Monochrom method), 44
 parse() (dipper.sources.MPD.MPD method), 42
 parse() (dipper.sources.MyChem.MyChem method), 44
 parse() (dipper.sources.MyDrug.MyDrug method), 45
 parse() (dipper.sources.NCBIGene.NCBIGene method), 46
 parse() (dipper.sources.OMIA.OMIA method), 47
 parse() (dipper.sources.OMIM.OMIM method), 48
 parse() (dipper.sources.Orphanet.Orphanet method), 49
 parse() (dipper.sources.OrthoXML.OrthoXML method), 49
 parse() (dipper.sources.Panther.Panther method), 50
 parse() (dipper.sources.Reactome.Reactome method), 52
 parse() (dipper.sources.RGD.RGD method), 51
 parse() (dipper.sources.SGD.SGD method), 52
 parse() (dipper.sources.Source.Source method), 54
 parse() (dipper.sources.StringDB.StringDB method), 55
 parse() (dipper.sources.UCSCBands.UCSCBands method), 56
 parse() (dipper.sources.UDP.UDP method), 57
 parse() (dipper.sources.WormBase.WormBase method), 58
 parse() (dipper.sources.ZFIN.ZFIN method), 59
 parse() (dipper.sources.ZFINSlim.ZFINSlim method), 60
 parse_checksum_file() (dipper.sources.IMPC.IMPC method), 39
 parse_mapping_file() (dipper.sources.Source.Source static method), 54
 Pathway (class in dipper.models.Pathway), 26
 pathway_parts (dipper.models.Pathway.Pathway attribute), 27
 PNTHDL (dipper.sources.Panther.Panther attribute), 50
 PostgreSQLSource (class in dipper.sources.PostgreSQLSource), 50
 prefix_exists() (dipper.utils.CurieUtil.CurieUtil method), 60
 process_all_common_disease_files() (dipper.sources.HPOAnnotations.HPOAnnotations method), 38
 process_allele_phenotype() (dipper.sources.WormBase.WormBase method), 58
 process_associations() (dipper.sources.OMIA.OMIA method), 47
 process_catalog() (dipper.sources.GWASCatalog.GWASCatalog method), 35
 process_classes() (dipper.sources.OMIA.OMIA method), 47
 process_common_disease_file() (dipper.sources.HPOAnnotations.HPOAnnotations method), 38
 process_disease_association() (dipper.sources.OMIA.OMIA method)

per.sources.WormBase.WormBase method),
 58
 process_entries() (dipper.sources.OMIM.OMIM
 method), 48
 process_feature_loc() (dip-
 per.sources.WormBase.WormBase method),
 58
 process_fish() (dipper.sources.ZFIN.ZFIN method), 59
 process_fish_disease_models() (dip-
 per.sources.ZFIN.ZFIN method), 60
 process_gaf() (dipper.sources.GeneOntology.GeneOntology
 method), 36
 process_gene_desc() (dip-
 per.sources.WormBase.WormBase
 58
 process_gene_ids() (dip-
 per.sources.WormBase.WormBase
 58
 process_gene_interaction() (dip-
 per.sources.WormBase.WormBase
 58
 process_mgi_note_allele_view() (dip-
 per.sources.MGI.MGI method), 40
 process_mgi_relationship_transgene_genes() (dip-
 per.sources.MGI.MGI method), 40
 process_nbk_html() (dip-
 per.sources.GeneReviews.GeneReviews
 method), 36
 process_omia_phenotypes() (dip-
 per.sources.Monarch.Monarch
 42
 process_orthology_evidence() (dip-
 per.sources.ZFIN.ZFIN method), 60
 process_pub_xrefs() (dip-
 per.sources.WormBase.WormBase
 58
 process_rnai_phenotypes() (dip-
 per.sources.WormBase.WormBase
 58
 process_species() (dipper.sources.OMIA.OMIA method),
 47
 process_xml_table() (dipper.sources.Source.Source
 method), 54
 properties (dipper.models.assoc.Association.Assoc
 attribute), 15
 properties (dipper.models.Environment.Environment
 attribute), 19
 properties (dipper.models.GenomicFeature.Feature
 attribute), 22
 properties (dipper.models.Genotype.Genotype attribute),
 24
 properties (dipper.models.Pathway.Pathway attribute), 27
 Provenance (class in dipper.models.Provenance), 27
 provenance_types (dip-
 per.models.Provenance.Provenance attribute),
 27

Q

querys (dipper.sources.FlyBase.FlyBase attribute), 34

R

RDFGraph (class in dipper.graph.RDFGraph), 13
 Reactome (class in dipper.sources.Reactome), 51
 REACTOME_BASE (dip-
 per.sources.Reactome.Reactome attribute),
 51
 ref_types (dipper.models.Reference.Reference attribute),
 28
 Reference (class in dipper.models.Reference), 27
 region_type_map (dipper.sources.Monochrom.Monochrom
 attribute), 44
 remove_backslash_r() (dipper.sources.Source.Source
 static method), 54
 remove_control_characters() (dip-
 per.utils.DipperUtil.DipperUtil
 method), 61
 replace() (in module dipper.utils.pysed), 62
 resources (dipper.sources.FlyBase.FlyBase attribute), 34
 resources (dipper.sources.MGI.MGI attribute), 40
 resources (dipper.sources.NCBIGene.NCBIGene attribute), 46
 return_target_list() (dipper.sources.MyChem.MyChem
 static method), 44
 RGD (class in dipper.sources.RGD), 51
 RGD_BASE (dipper.sources.RGD.RGD attribute), 51
 rmlinematch() (in module dipper.utils.pysed), 62
 rmlinenumber() (in module dipper.utils.pysed), 62
 RomanError, 62

S

SCIGRAPH_BASE (dip-
 per.sources.NCBIGene.NCBIGene attribute),
 45
 scrub() (dipper.sources.ClinVar.ClinVar method), 31
 scrub() (dipper.sources.HPOAnnotations.HPOAnnotations
 method), 38
 scrub() (dipper.sources.OMIA.OMIA method), 47
 scrub() (dipper.sources.ZFIN.ZFIN method), 60
 serialize() (dipper.graph.Graph.Graph method), 13
 serialize() (dipper.graph.StreamedGraph.StreamedGraph
 method), 14
 set_association_id() (dip-
 per.models.assoc.Association.Assoc
 method), 15
 set_association_id() (dip-
 per.models.assoc.Chem2DiseaseAssoc.Chem2DiseaseAssoc
 method), 15

set_association_id()	(dip-	Source (class in dipper.sources.Source), 52
per.models.assoc.D2PAssoc.D2PAssoc method), 16		species (dipper.sources.WormBase.WormBase attribute), 58
set_association_id()	(dip-	static_files (dipper.sources.CTD.CTD attribute), 31
per.models.assoc.G2PAssoc.G2PAssoc method), 17		StreamedGraph (class in dipper.graph.StreamedGraph), 14
set_citation() (dipper.models.Dataset.Dataset method), 18		STRING_BASE (dipper.sources.StringDB.StringDB at- tribute), 55
set_date_issued() (dipper.models.Dataset.Dataset method), 18		StringDB (class in dipper.sources.StringDB), 55
set_description() (dipper.models.assoc.Association.Assoc method), 15		
set_environment()	(dip-	T
per.models.assoc.G2PAssoc.G2PAssoc method), 17		tables (dipper.sources.EOM.EOM attribute), 33
set_license() (dipper.models.Dataset.Dataset method), 18		tables (dipper.sources.FlyBase.FlyBase attribute), 34
set_object() (dipper.models.assoc.Association.Assoc method), 15		terms (dipper.models.assoc.OntologyAssoc.OntologyAssoc attribute), 18
set_relationship() (dipper.models.assoc.Association.Assoc method), 15		terms (dipper.sources.Coriell.Coriell attribute), 32
set_score() (dipper.models.assoc.Association.Assoc method), 15		terms (dipper.sources.GWASCatalog.GWASCatalog at- tribute), 35
set_stage() (dipper.models.assoc.G2PAssoc.G2PAssoc method), 17		test_graph_equality() (dipper.utils.TestTools.TestTools method), 62
set_subject() (dipper.models.assoc.Association.Assoc method), 15		test_ids (dipper.sources.AnimalQTLdb.AnimalQTLdb at- tribute), 29
set_version_by_date() (dipper.models.Dataset.Dataset method), 18		test_ids (dipper.sources.IMPC.IMPC attribute), 39
set_version_by_num() (dipper.models.Dataset.Dataset method), 18		test_ids (dipper.sources.KEGG.KEGG attribute), 40
setAuthorList() (dipper.models.Reference.Reference method), 28		test_ids (dipper.sources.MMRRC.MMRRC attribute), 41
setFileAccessUrl() (dipper.models.Dataset.Dataset method), 18		test_ids (dipper.sources.MPD.MPD attribute), 42
setShortCitation() (dipper.models.Reference.Reference method), 28		test_ids (dipper.sources.OMIM.OMIM attribute), 48
settestmode() (dipper.sources.Source.Source method), 54		test_ids (dipper.sources.WormBase.WormBase attribute), 58
settestonly() (dipper.sources.Source.Source method), 55		test_ids (dipper.sources.ZFIN.ZFIN attribute), 60
setTitle() (dipper.models.Reference.Reference method), 28		test_keys (dipper.sources.FlyBase.FlyBase attribute), 34
setType() (dipper.models.Reference.Reference method), 28		test_keys (dipper.sources.MGI.MGI attribute), 40
setVersion() (dipper.models.Dataset.Dataset method), 18		test_lines (dipper.sources.Coriell.Coriell attribute), 32
setYear() (dipper.models.Reference.Reference method), 28		TestUtils (class in dipper.utils.TestTools), 62
SGD (class in dipper.sources.SGD), 52		toRoman() (in module dipper.utils.romanplus), 63
SGD_BASE (dipper.sources.SGD.SGD attribute), 52		types (dipper.models.GenomicFeature.Feature attribute), 22
skolemizeBlankNode() (dipper.graph.Graph.Graph method), 13		types (dipper.models.Model.Model attribute), 26
skolemizeBlankNode()	(dip-	
per.graph.RDFGraph.RDFGraph	method),	U
14		UCSCBands (class in dipper.sources.UCSCBands), 55
skolemizeBlankNode()	(dip-	UDP (class in dipper.sources.UDP), 56
per.graph.StreamedGraph.StreamedGraph method), 14		UDP_SERVER (dipper.sources.UDP.UDP attribute), 57
		update_wsnum_in_files() (dip- per.sources.WormBase.WormBase method), 58
		V
		variant_ids (dipper.sources.ClinVar.ClinVar attribute), 31
		W
		wbdev (dipper.sources.WormBase.WormBase attribute), 59

wbprod (dipper.sources.WormBase.WormBase attribute),
 59
wbre1 (dipper.sources.WormBase.WormBase attribute),
 59
whoami() (dipper.sources.Source method), 55
WormBase (class in dipper.sources.WormBase), 57
write() (dipper.sources.Source method), 55
write() (dipper.utils.GraphUtils.GraphUtils method), 62
write_molgen_report() (dipper.sources.OMIA.OMIA
method), 47

Z

ZFIN (class in dipper.sources.ZFIN), 59
ZFINSlim (class in dipper.sources.ZFINSlim), 60
zygosity (dipper.models.Genotype.Genotype attribute),
 25