

---

# Dipper Documentation

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

**Dec 27, 2021**



---

## Contents

---

<b>1 Getting started</b>	<b>3</b>
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
<b>2 Deeper into Dipper</b>	<b>7</b>
2.1 Working with graphs . . . . .	7
2.2 Working with the model 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
<b>3 For developers</b>	<b>13</b>
3.1 API Docs . . . . .	13
3.2 Source APIs . . . . .	65
<b>4 Indices and tables</b>	<b>67</b>
<b>Python Module Index</b>	<b>69</b>
<b>Index</b>	<b>71</b>



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 impc,hpoa
```

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 model 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 model 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 the `serialize` method inherited from the parent `RDFLib` `graph` class<sup>3</sup>:

```
from dipper.graph.RDFGraph import RDFGraph

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

(continues on next page)

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

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

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

(continued from previous page)

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

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 model API

The model 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 model 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 Reification<sup>1</sup> pattern to create ternary statements, for example, adding frequency data to phenotype to disease associations. We utilize the Open Biomedical Association ontology<sup>2</sup> 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 ontology<sup>4</sup> 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'
    },
    ...
}
```

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

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

<sup>4</sup> 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.yaml` 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”: “myomimkey1234”, “omimftp”: “myomimftpkey5678”, “ncbi”: “myncbikey”,  
}, “dbauth”: {  
    “mgi”: { ‘user’: ‘mymgiuser’, ‘password’: ‘mymgipw’, ‘host’: ‘mgi-adhoc.jax.org’,  
        ‘database’: ‘mgd’, ‘port’: 5432  
    }  
}  
}
```

This file must be placed in the dipper package directory and named `conf.yaml`. 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
        Graph class, used by RDFGraph and StreamedGraph
        addTriple(subject_id, predicate_id, obj, object_is_literal=False, literal_type=None, subject_category=None, object_category=None)
        curie_regex = re.compile('^[a-zA-Z_]?[a-zA-Z_0-9_]*:[A-Za-z0-9_.-]*[A-Za-z0-9_.-]*')
        serialize(**kwargs)
        skolemizeBlankNode(curie)
```

##### dipper.graph.RDFGraph module

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

---

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=None*, *literal\_type=None*, *subject\_category=None*, *object\_category=None*)

**bind\_all\_namespaces()**

Results in the RDF @prefix directives for every ingest being added to this ingest.

**curie\_map** = { ''': 'https://monarchinitiative.org/' , 'APB': 'http://pb.apf.edu.au/phenb'}

**curie\_util** = <*dipper.utils.CurieUtil.CurieUtil* object>

**fhandle** = <*\_io.TextIOWrapper* name='/home/docs/checkouts/readthedocs.org/user\_builds/dip'>

**globaltcid** = {':activating': 'activating\_mutation', ':all\_missense\_or\_inframe': 'all'}

**globaltt** = {'3\_prime\_UTR\_variant': 'SO:0001624', '3prime\_overlapping\_ncRNA': 'SO:0002'}

**serialize** (*destination=None*, *format='turtle'*, *base=None*, *encoding=None*)

Serialize the Graph to destination

If destination is None serialize method returns the serialization as bytes or string.

If encoding is None and destination is None, returns a string If encoding is set, and Destination is None, returns bytes

Format defaults to turtle.

Format support can be extended with plugins, but “xml”, “n3”, “turtle”, “nt”, “pretty-xml”, “trix”, “trig” and “nquads” are built in.

**skolemizeBlankNode** (*curie*)

## dipper.graph.StreamedGraph module

**class** *dipper.graph.StreamedGraph.StreamedGraph* (*are\_bnodes\_skized=True*, *identi-*  
*fier=None*, *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*, *obj*, *object\_is\_literal=None*, *literal\_type=None*, *sub-*  
*ject\_category=None*, *object\_category=None*)

**curie\_map** = { ''': 'https://monarchinitiative.org/' , 'APB': 'http://pb.apf.edu.au/phenb'}

**curie\_util** = <*dipper.utils.CurieUtil.CurieUtil* object>

**fhandle** = <*\_io.TextIOWrapper* name='/home/docs/checkouts/readthedocs.org/user\_builds/dip'>

**globaltcid** = {':activating': 'activating\_mutation', ':all\_missense\_or\_inframe': 'all'}

**globaltt** = {'3\_prime\_UTR\_variant': 'SO:0001624', '3prime\_overlapping\_ncRNA': 'SO:0002'}

**serialize** (*subject\_iri*, *predicate\_iri*, *obj*, *object\_is\_literal=False*, *literal\_type=None*,  
*subject\_category\_iri=None*, *predicate\_category\_iri='biolink:category'*, *ob-*  
*ject\_category\_iri=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, subject_category=None, object_category=None)
```

Bases: object

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

**add\_association\_to\_graph**(*association\_category*=None)

**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**

**get\_association\_id**()

**static make\_association\_id**(*definedby*, *sub*, *pred*, *obj*, *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.

Note this is equivalent to a RDF blank node

**Parameters**

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

**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\_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, definedby, 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.D2PAssoc(graph, definedby, disease\_id, phenotype\_id, onset=None, frequency=None, rel=None, disease\_category=None, phenotype\_category=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**(association\_category=None)

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 –

- disease\_category – a biolink category CURIE for disease\_id (defaults to

biolink:Disease via the constructor) :param phenotype\_category: a biolink category CURIE for phenotype\_id (defaults to biolink:PhenotypicFeature via the constructor) :return:

### **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.G2PAssoc module

```
class dipper.models.assoc.G2PAssoc(graph, definedby, entity_id, phenotype_id,
                                     rel=None, entity_category=None, phenotype_category=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(entity\_category=None, phenotype\_category=None)**

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: :param entity\_category: a biolink category CURIE for self.sub :param phenotype\_category: a biolink category CURIE for self.obj :return:

### **make\_g2p\_id()**

Make an association id for phenotypic associations that is defined by: source of association + (Annot 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
```

### dipper.models.assoc.OrthologyAssoc module

```
class dipper.models.assoc.OrthologyAssoc.OrthologyAssoc(graph, definedby, gene1,
                                                               gene2, rel=None, subject_category=None,
                                                               object_category=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 EDAM-DATA:gene\_family <family\_id> RO:has\_member <gene1> <family\_id> RO:has\_member <gene2> <gene1> biolink:category <subject\_category> <gene2> biolink:category <object\_category> :param family\_id: :param g: the graph to modify :return:

## Submodules

### dipper.models.BiolinkVocabulary module

```
class dipper.models.BiolinkVocabulary.BioLinkVocabulary
Bases: object
```

```
bl_file_with_path = '/home/docs/checkouts/readthedocs.org/user_builds/dipper/checkouts/
bl_vocab = {'curie_prefix': 'biolink', 'terms': ['AnatomicalEntity', 'Association',
key = 'Zygosity'
terms = {'AnatomicalEntity': 'biolink:AnatomicalEntity', 'Association': 'biolink:Ass
yaml_file = <_io.TextIOWrapper name='/home/docs/checkouts/readthedocs.org/user_builds/
```

### dipper.models.ClinVarRecord module

<https://www.ncbi.nlm.nih.gov/clinvar/docs/details/> Object mapping to XML schema

```
class dipper.models.ClinVarRecord.Allele(id: str, label: Optional[str] = None, variant_type: Optional[str] = None, genes: Optional[List[dipper.models.ClinVarRecord.Gene]] = None, synonyms: Optional[List[str]] = None, dbsnps: Optional[List[str]] = None)
Bases: object
```

ClinVar Allele Alleles can have 0 to many genes

These are called alleles and variants on the ClinVar UI, and variant, single nucleotide variant, etc in the XML  
id: allele id label: label variant\_type: single nucleotide variant genes: gene(s) in which the variant is found  
synonyms: eg HGVC dbsnp: dbSNP curies

---

```
class dipper.models.ClinVarRecord.ClinVarRecord(id: str, accession: str, created: str, updated: str, genovar: dipper.models.ClinVarRecord.Genovar, significance: str, conditions: Optional[List[dipper.models.ClinVarRecord.Condition]] = None)
```

Bases: object

Reference ClinVar Record (RCV) id: RCV id accession: RCV accession (eg RCV000123456) created: Created date updated: Updated date genovar: the variant or genotype associated with the condition(s) significance: clinical significance (eg benign, pathogenic) condition: The condition(s) for which this allele set was interpreted, with

links to databases with defining information about that condition.

```
class dipper.models.ClinVarRecord.Condition(id: str, label: Optional[str] = None, database: Optional[str] = None, medgen_id: Optional[str] = None)
```

Bases: object

ClinVar condition

```
class dipper.models.ClinVarRecord.Gene(id: Union[str, int, None], association_to_allele: str)
```

Bases: object

ClinVar Gene Intentionally leaves out label/symbol, this should come from HGNC

```
class dipper.models.ClinVarRecord.Genotype(id: str, label: Optional[str] = None, variants: Optional[List[dipper.models.ClinVarRecord.Variant]] = None, variant_type: Optional[str] = None)
```

Bases: *dipper.models.ClinVarRecord.Genovar*

ClinVar genotype Example: Compound Heterozygote, Diplootype

These are called variants on the ClinVar UI, and a GenotypeSet in the XML

```
class dipper.models.ClinVarRecord.Genovar(id: str, label: Optional[str] = None, variant_type: Optional[str] = None)
```

Bases: object

Sequence feature entity that is linked to a disease in an Reference ClinVar Record, it is either a Variant or Genotype

```
class dipper.models.ClinVarRecord.Variant(id: str, label: Optional[str] = None, alleles: Optional[List[dipper.models.ClinVarRecord.Allele]] = None, variant_type: Optional[str] = None)
```

Bases: *dipper.models.ClinVarRecord.Genovar*

ClinVar variant, variants can have one or more alleles

These are called variants and alleles on the ClinVar UI, and Variants in the XML

## dipper.models.Dataset module

Produces metadata about ingested data

```
class dipper.models.Dataset.Dataset(identifier,      data_release_version,      ingest_name,
                                    ingest_title,      ingest_url,      ingest_logo=None,      in-
                                    gest_description=None,      license_url=None,
                                    data_rights=None,      graph_type='rdf_graph',
                                    file_handle=None,      distribution_type='ttl',
                                    dataset_curie_prefix='MonarchArchive')
```

Bases: object

This class produces metadata about a dataset that is compliant with the HCLS dataset specification: [https://www.w3.org/TR/2015/NOTE-hcls-dataset-20150514/#s4\\_4](https://www.w3.org/TR/2015/NOTE-hcls-dataset-20150514/#s4_4)

Summary level: The summary level provides a description of a dataset that is independent of a specific version or format. (e.g. the Monarch ingest of CTD) CURIE for this is something like MonarchData:[SOURCE IDENTIFIER]

Version level: The version level captures version-specific characteristics of a dataset. (e.g. the 01-02-2018 ingest of CTD) CURIE for this is something like MonarchData:[SOURCE IDENTIFIER\_INGESTTIMESTAMP]

Distribution level: The distribution level captures metadata about a specific form and version of a dataset (e.g. turtle file for 01-02-2018 ingest of CTD). There is a [distribution level resource] for each different downloadable file we emit, i.e. one for the TTL file, one for the ntriples file, etc. CURIE for this is like MonarchData:[SOURCE IDENTIFIER\_INGESTTIMESTAMP].ttl or MonarchData:[SOURCE IDENTIFIER\_INGESTTIMESTAMP].nt or MonarchData:[SOURCE IDENTIFIER\_INGESTTIMESTAMP].[whatever file format]

We write out at least the following triples:

SUMMARY LEVEL TRIPLES: [summary level resource] - rdf:type -> dcterms:Dataset [summary level resource] - dc:title -> title (literal) [summary level resource] - dc:description -> description (literal)

(use docstring from Source class)

[summary level resource] - dc:source -> [source web page, e.g. omim.org] [summary level resource] - schema:logo -> [source logo IRI] [summary level resource] - dc:publisher -> monarchinitiative.org

n.b: about summary level resource triples: – HCLS spec says we “should” link to our logo and web page, but I’m not, because it would confuse the issue of whether we are pointing to our logo/page or the logo/page of the data source for this ingest. Same below for [version level resource] and [distibution level resource] - I’m not linking to our page/logo down there either. – spec says we “should” include summary level triples describing Update frequency and SPARQL endpoint but I’m omitting this for now, because these are not clearly defined at the moment

VERSION LEVEL TRIPLES: [version level resource] - rdf:type -> dcterms:Dataset [version level resource] - dc:title -> version title (literal) [version level resource] - dc:description -> version description (literal) [version level resource] - dc:created -> ingest timestamp [ISO 8601 compliant] [version level resource] - pav:version -> ingest timestamp (same one above) [version level resource] - dc:creator -> monarchinitiative.org [version level resource] - dc:publisher -> monarchinitiative.org [version level resource] - dc:isVersionOf -> [summary level resource] [version level resource] - dc:source -> [source file 1 IRI] [version level resource] - dc:source -> [source file 2 IRI] ...

[source file 1 IRI] - pav:retrievedOn -> [download date timestamp] [source file 2 IRI] - pav:version -> [source version (if set, optional)] [source file 2 IRI] - pav:retrievedOn -> [download date timestamp] [source file 2 IRI] - pav:version -> [source version (if set, optional)] ...

[version level resource] - pav:createdWith -> [Dipper github URI] [version level resource] - void:dataset -> [distribution level resource]

[version level resource] - cito:citesAsAuthoriy -> [citation id 1] [version level resource] - cito:citesAsAuthoriy -> [citation id 2] [version level resource] - cito:citesAsAuthoriy -> [citation id 3]

n.b: about version level resource triples: - spec says we “should” include Date of issue/dc:issued triple, but I’m not because it is redundant with this triple above: [version level resource] - dc:created -> time stamp and would introduce ambiguity and confusion if the two disagree. Same below for [distribution level resource] - dc:created -> tgiime stamp below Also omitting:

- triples linking to our logo and page, see above.
- License/dc:license triple, because we will make this triple via the [distribution level resource] below
- Language/dc:language triple b/c it seems superfluous. Same below for [distribution level resource] - no language triple.
- [version level resource] - pav:version triple is also a bit redundant

with the pav:version triple below, but the spec requires both these triples - I’m omitting the [version level resource] -> pav:previousVersion because Dipper doesn’t know this info for certain at run time. Same below for [distribution level resource] - pav:previousVersion.

DISTRIBUTION LEVEL TRIPLES: [distribution level resource] - rdf:type -> dctypes:Dataset [distribution level resource] - rdf:type -> dcat:Distribution [distribution level resource] - dc:title -> distribution title (literal) [distribution level resource] - dc:description -> distribution description (lit.) [distribution level resource] - dc:created -> ingest timestamp[ISO 8601 compliant] [distribution level resource] - pav:version -> ingest timestamp (same as above) [distribution level resource] - dc:creator -> monarchinitiative.org [distribution level resource] - dc:publisher -> monarchinitiative.org [distribution level resource] - dc:license -> [license info, if available

otherwise indicate unknown]

[distribution level resource] - dc:rights -> [data rights IRI] [distribution level resource] - pav:createdWith -> [Dipper github URI] [distribution level resource] - dc:format -> [IRI of ttl|nt|whatever spec] [distribution level resource] - dcat:downloadURL -> [ttl|nt URI] [distribution level resource] - void:triples -> [triples count (literal)] [distribution level resource] - void:entities -> [entities count (literal)] [distribution level resource] - void:distinctSubjects -> [subject count (literal)] [distribution level resource] - void:distinctObjects -> [object count (literal)] [distribution level resource] - void:properties -> [properties count (literal)] ...

n.b: about distribution level resource triples: - omitting Vocabularies used/void:vocabulary and Standards used/dc:conformTo triples, because they are described in the ttl file - also omitting Example identifier/idot:exampleIdentifier and Example resource/void:exampleResource, because we don’t really have one canonical example of either - they’re all very different. - [distribution level resource] - dc:created should have the exact same time stamp as this triple above: [version level resource] - dc:created -> time stamp - this [distribution level resource] - pav:version triple should have the same object as [version level resource] - pav:version triple above - Data source provenance/dc:source triples are above in the [version level resource] - omitting Byte size/dc:byteSize, RDF File URL/void:dataDump, and Linkset/void:subset triples because they probably aren’t necessary for MI right now - these triples “should” be emitted, but we will do this in a later iteration: # of classes void:classPartition IRI # of literals void:classPartition IRI # of RDF graphs void:classPartition IRI

Note: Do not use blank nodes in the dataset graph. This dataset graph is added to the main Dipper graph in Source.write() like so

```
$ mainGraph = mainGraph + datasetGraph
```

which apparently in theory could lead to blank node ID collisions between the two graphs.

Note also that this implementation currently does not support producing metadata for StreamedGraph graphs (see dipper/graph/StreamedGraph.py). StreamedGraph is currently not being used for any ingest, so this isn’t

a problem. There was talk of using StreamedGraph for a rewrite/refactor of the Clinvar ingest, which would probably require adding support here for StreamedGraph's.

### `get_graph()`

This method returns the dataset graph :param :return: dataset graph

### `get_license()`

This method returns the license info :param :return: license info

### `static hash_id(word)`

Given a string, make a hash Duplicated from Source.py.

**Parameters** `word` – str string to be hashed

**Returns** 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 Duplicated from Source.py to avoid circular imports. :param long\_string: string to use to generate identifier :param prefix: prefix to prepend to identifier [Monarch] :return: a Monarch identifier

### `set_citation(citation_id)`

This method adds [citaton\_id] argument to the set of citations, and also adds a triple indicating that version level cito:citesAsAuthority [citation\_id] :param: citation\_id :return: none

### `set_ingest_source(url, predicate=None, is_object_literal=False)`

This method writes a triple to the dataset graph indicating that the ingest used a file or resource at [url] during the ingest.

Triple emitted is version\_level\_curie dc:source [url]

This triple is likely to be redundant if Source.get\_files() is used to retrieve the remote files/resources, since this triple should also be emitted as files/resources are being retrieved. This method is provided as a convenience method for sources that do their own downloading of files.

**Parameters**

- `url` – a remote resource used as a source during ingest
- `predicate` – the predicate to use for the triple ["dc:source"] from spec (<https://www.w3.org/TR/2015/NOTE-hcls-dataset-20150514/>) “Use dc:source when the source dataset was used in whole or in part. Use pav:retrievedFrom when the source dataset was used in whole and was not modified from its original distribution. Use prov:wasDerivedFrom when the source dataset was in whole or in part and was modified from its original distribution.”

**Returns** None

### `set_ingest_source_file_version_date(file_iri, date, datatype=rdflib.term.URIRef('http://www.w3.org/2001/XMLSchema#date'))`

This method sets the version that the source (OMIM, CTD, whatever) uses to refer to this version of the remote file/resource that was used in the ingest

It writes this triple:

`file_iri - 'pav:version' -> date or timestamp`

Version is added as a literal of datatype XSD date

Note: if `file_iri` was retrieved using `get_files()`, then the following triple was created and you might not need this method:

`file_iri - 'pav:retrievedOn' -> download date`

**Parameters**

- **file\_iri** – a remote file or resource used in ingest
- **date** – a date in YYYYMMDD format that the source (OMIM, CTD). You can add timestamp as a version by using a different datatype (below) :param datatype: an XSD literal datatype, default is XSD.date uses to refer to this version of the file/resource used during the ingest :return: None

**set\_ingest\_source\_file\_version\_num(file\_iri, version)**

This method sets the version of a remote file or resource that is used in the ingest. It writes this triple:

file\_iri - ‘pav:version’ -> version

Version is an untyped literal

Note: if your version is a date or timestamp, use set\_ingest\_source\_file\_version\_date() instead

**Parameters**

- **file\_iri** – a remote file or resource used in ingest
- **version** – a number or string (e.g. v1.2.3) that the source (OMIM, CTD)

uses to refer to this version of the file/resource used during the ingest :return: None

**set\_ingest\_source\_file\_version\_retrieved\_on(file\_iri,**

*date,*  
*datatype=rdf:term.URIRef('http://www.w3.org/2001/XMLSchema#dateTime')*

This method sets the date on which a remote file/resource (from OMIM, CTD, etc) was retrieved.

It writes this triple:

file\_iri - ‘pav:retrievedOn’ -> date or timestamp

Version is added as a literal of datatype XSD date by default

Note: if file\_iri was retrieved using get\_files(), then the following triple was created and you might not need this method:

file\_iri - ‘pav:retrievedOn’ -> download date

**Parameters**

- **file\_iri** – a remote file or resource used in ingest
- **date** – a date in YYYYMMDD format that the source (OMIM, CTD). You can

add timestamp as a version by using a different datatype (below) :param datatype: an XSD literal datatype, default is XSD.date uses to refer to this version of the file/resource used during the ingest :return: None

## 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 its parts to a supplied graph.

This is a stub.

**addComponentAttributes(component\_id, entity\_id, value=None, unit=None, component\_category=None, entity\_category=None)**

**addComponentToEnvironment(env\_id, component\_id, environment\_category=None, component\_category=None)**

**addEnvironment(env\_id, env\_label, env\_type=None, env\_description=None)**

```
addEnvironmentalCondition(cond_id, cond_label, cond_type=None, cond_description=None, condition_category=None)
```

### dipper.models.Evidence module

```
class dipper.models.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, data_curie_category=None)  
Add data individual :param data_curie: str either curie formatted or long string,
```

long strings will be converted to bnodes

#### Parameters

- **data\_curie\_category** – a biolink category CURIE for data\_curie
- **type** – str curie
- **label** – str

#### Returns

None

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

#### Parameters

- **evidence\_line** – curie or iri, evidence line
- **evidence\_type** – curie or iri, evidence type if available

#### Returns

None

```
add_source(evidence_line, source, label=None, src_type=None, source_category=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, evidence_type=None, label=None)  
Add supporting line of evidence node to association id
```

**Parameters**

- **evidence\_line** – curie or iri, evidence line
- **evidence\_type** – curie or iri, evidence type if available

**Returns** None

```
add_supporting_publication(evidence_line, publication, label=None, pub_type=None)
<evidence> <has_supporting_reference> <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:
```

**dipper.models.Family module**

```
class dipper.models.Family.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, group_category=None, member_category=None)
addMemberOf(member_id, group_id, member_category=None, group_category=None)
```

**dipper.models.GenomicFeature module**

```
class dipper.models.GenomicFeature.Feature(graph, feature_id=None, label=None, fea-
ture_type=None, description=None, fea-
ture_category=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:
```

```
addFeatureProperty(property_type, feature_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:
```

```
addFeatureToGraph(add_region=True, region_id=None, feature_as_class=False, fea-
ture_category=None)
```

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
```

```
:param add_region [True] :param region_id [None] :param feature_as_class [False] :param feature_category: a biolink category CURIE for feature
```

### **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*, *subject\_category=None*, *object\_category=None*)

This will add reciprocal triples like: feature <is subsequence of> parent parent has\_subsequence feature :param graph: :param parentid:

#### Parameters

### **addTaxonToFeature** (*taxonid*)

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

### 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:has\_affected\_feature.

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\_allele\_of.

Here, the allele should be a variant\_locus, not a sequence alteration. :param allele\_id: :param gene\_id: :param rel\_id: :return:

**addChromosome** (*chrom*, *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*, *construct\_category=None*, *construct\_type\_category=None*)**Parameters**

- **construct\_id** –
- **construct\_label** –
- **construct\_type** –
- **construct\_description** –
- **construct\_category** – a biolink category CURIE for construct\_id
- **construct\_type\_category** – a biolink category CURIE for construct\_type

**Returns****addDerivesFrom** (*child\_id*, *parent\_id*, *child\_category=None*, *parent\_category=None*)

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 its 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=None*, *gene\_type=None*, *gene\_description=None*)  
genes are classes**addGeneProduct** (*sequence\_id*, *product\_id*, *product\_label=None*, *product\_type=None*, *sequence\_category=None*, *product\_category=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: :param sequence\_category: bl category CURIE for seq\_id [blv.terms.Gene].value :param product\_category: biolink category CURIE for product\_id :return:

**addGeneTargetingReagent** (*reagent\_id*, *reagent\_label*, *reagent\_type*, *gene\_id*, *description=None*, *reagent\_category=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*)

Add genotype has\_variant\_part reagent\_id. For example, add a morphant reagent thingy to the genotype, assuming it's a extrinsic\_genotype Also a triple to assign biolink categories to genotype and reagent. :param reagent\_id :param genotype\_id :return:

**addGenome** (*taxon\_num*, *taxon\_label=None*, *genome\_id=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*, *part\_category=None*, *parent\_category=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: :param part\_category: a biolink vocab curie for part\_id :param parent\_category: a biolink vocab curie for parent\_id :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*, *targeted\_gene\_label=None*, *description=None*, *reagent\_category=None*)

This will create the instance of a gene that is targeted by a molecular reagent (such as a morpholino or rna). 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 rdfs:label targeted\_gene\_label dc:description description <reagent\_id> GENO:targets\_gene <gene\_id>

### Parameters

- **reagent\_id** –
- **gene\_id** –
- **targeted\_gene\_id** –
- **reagent\_category** – a biolink category CURIE for reagent\_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, child_category=None, parent_category=None)
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, genopart_category=None)
```

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: :param genopart\_category: a biolink term for genopart\_id :return:

```
addVSLCtoParent (vslc_id, parent_id, part_category=None, parent_category=None)
```

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: :param part\_category: a biolink category CURIE for part :param parent\_category: a biolink category CURIE for parent :return:

```
static makeGenomeID (taxon_id)
```

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

```
static 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:

## dipper.models.Model module

```
class dipper.models.Model.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, class_category=None, class_type_category=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** –
- **class\_category** – a biolink category CURIE for class
- **class\_type\_category** – a biolink category CURIE for class type

### Returns

**addComment** (*subject\_id*, *comment*, *subject\_category=None*)

**addDefinition** (*class\_id*, *definition*, *class\_category=None*)

**addDepiction** (*subject\_id*, *image\_url*)

**addDeprecatedClass** (*old\_id*, *new\_ids=None*, *old\_id\_category=None*, *new\_ids\_category=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.

:param old\_id\_category - a biolink category CURIE for old id :param new\_ids\_category - a biolink category CURIE for new ids :return: None

**addDeprecatedIndividual** (*old\_id*, *new\_ids=None*, *old\_id\_category=None*, *new\_id\_category=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.

:param old\_id\_category - a biolink category CURIE for old id :param new\_ids\_category - a biolink category CURIE for new ids :return:

**addDescription** (*subject\_id*, *description*, *subject\_category=None*)

**addEquivalentClass** (*sub*, *obj*, *subject\_category=None*, *object\_category=None*)

**addIndividualToGraph** (*ind\_id*, *label*, *ind\_type=None*, *description=None*, *ind\_category=None*, *ind\_type\_category=None*)

**addLabel** (*subject\_id*, *label*, *subject\_category=None*)

**addOWLPropertyClassRestriction** (*class\_id*, *property\_id*, *property\_value*, *class\_category=None*, *property\_id\_category=None*, *property\_value\_category=None*)

**addOWLVersionIRI** (*ontology\_id*, *version\_iri*)

**addOWLVersionInfo** (*ontology\_id*, *version\_info*)

**addOntologyDeclaration** (*ontology\_id*)

---

**addPerson** (*person\_id*, *person\_label=None*)  
**addSameIndividual** (*sub*, *obj*, *subject\_category=None*, *object\_category=None*)  
**addSubClass** (*child\_id*, *parent\_id*, *child\_category=None*, *parent\_category=None*)  
**addSynonym** (*class\_id*, *synonym*, *synonym\_type=None*, *class\_category=None*)  
 Add the synonym as a property of the class cid. Assume it is an exact synonym, unless otherwise specified  
 :param self: :param class\_id: class id :param synonym: the literal synonym label :param synonym\_type:  
 the CURIE of the synonym type (not the URI) :param class\_category: biolink category CURIE for class\_id  
 (no biolink category is possible for synonym, since this is added to the triple as a literal) :return:  
**addTriple** (*subject\_id*, *predicate\_id*, *obj*, *object\_is\_literal=False*, *literal\_type=None*, *sub-*  
*ject\_category=None*, *object\_category=None*)  
**addType** (*subject\_id*, *subject\_type*, *subject\_category=None*, *subject\_type\_category=None*)  
**addXref** (*class\_id*, *xref\_id*, *xref\_as\_literal=False*, *class\_category=None*, *xref\_category=None*)  
**makeLeader** (*node\_id*)  
 Add an annotation property to the given `node\_id` to be the clique\_leader. This is a monarchism.  
 :param node\_id: :param node\_category: a biolink category CURIE for node\_id :return:

## dipper.models.Pathway module

**class** dipper.models.Pathway.**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** –
- **component\_category** – biolink category for component\_id
- **pathway\_category** – biolink category for pathway\_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:

### 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,
                    agent_category=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, object_is_literal=None)
add_study_parts(study, study_parts, study_parts_category=None)
add_study_to_measurements(study, measurements)
```

### 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, subject_category=None, page_category=None)
```

```
addRefToGraph()
```

```
addTitle(subject_id, title)
```

```
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.AnimalQTLdb(graph_type, are_bnodes_skolemized,
                                                data_release_version=None)
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.

```
GENEINFO = 'ftp://ftp.ncbi.nih.gov/gene/DATA/GENE_INFO'
GITDIP = 'https://raw.githubusercontent.com/monarch-initiative/dipper/master'
fetch(is_dl_forced=False)
    abstract method to fetch all data from an external resource. this should be overridden by subclasses :return:
    None

files = {'Bos_taurus_info': {'columns': ['tax_id', 'GeneID', 'Symbol', 'LocusTag', ''],
                             'gene_info_columns': ['tax_id', 'GeneID', 'Symbol', 'LocusTag', 'Synonyms', 'dbXrefs'],
                             'getTestSuite': ()}
        }
    An abstract method that should be overwritten with tests appropriate for the specific source. :return:
    gff_columns = ['SEQNAME', 'SOURCE', 'FEATURE', 'START', 'END', 'SCORE', 'STRAND', 'FRAMES']
    parse(limit=None)

    Parameters limit -
    Returns
    qtl_columns = ['QTL_ID', 'QTL_symbol', 'Trait_name', 'assotype', '(empty)', 'Chromosome',
                   'test_ids = {1795, 1798, 8945, 12532, 14234, 17138, 28483, 29016, 29018, 29385, 31023, }
```

```
trait_mapping_columns = ['VT', 'LPT', 'CMO', 'ATO', 'Species', 'Class', 'Type', 'QTL_C
```

### dipper.sources.Bgee module

```
class dipper.sources.Bgee(graph_type, are_bnodes_skolemized,
                           data_release_version=None, 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.

```
check_if_remote_is_newer(localfile, remote_size, remote_modify)
Overrides check_if_remote_is_newer 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

```
default_species = ['Cavia porcellus', 'Mus musculus', 'Rattus norvegicus', 'Monodelphis
fetch(is_dl_forced=False)
```

**Parameters** **is\_dl\_forced** – boolean, force download

#### Returns

```
files = {'anat_entity': {'columns': ['Ensembl gene ID', 'gene name', 'anatomical ent
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(graph_type, are_bnodes_skolemized,
                           data_release_version=None, 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': 'BIOGRID-IDENTIFIERS-LATEST.tab.zip', 'url': 'http://
```

---

```
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(graph_type, are_bnodes_skolemized, data_release_version=None)
Bases: dipper.sources.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. (updated monthly).

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’ or ‘therapeutic’. (observe strictly OR) Unfortunately, we cannot disambiguate between marker (gene expression) and mechanism (causation) for these associations. Therefore, we are left to relate these simply by “marker”.

# We DISCONTINUED at some point prior to 202005 # 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_associations': {'columns': ['ChemicalName', 'ChemicalID',
```

```
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
```

## dipper.sources.ClinVar module

Converts ClinVar XML into RDF triples to be ingested by SciGraph. These triples conform to the core of the SEPIO Evidence & Provenance model

We also use the clinvar curated gene to disease mappings to discern the functional consequence of a variant on a gene in cases where this is ambiguous. For example, some variants are located in two genes overlapping on different

strands, and may only have a functional consequence on one gene. This is suboptimal and we should look for a source that directly provides this.

**creating a test set.** get a full dataset default ClinVarFullRelease\_00-latest.xml.gz get the mapping file default gene\_condition\_source\_id get a list of RCV default CV\_test\_RCV.txt put the input files the raw directory write the test set back to the raw directory

```
./scripts/ClinVarXML_Subset.sh | gzip > raw/clinvar/ClinVarTestSet.xml.gz
```

parsing a test set (Skolemizing blank nodes i.e. for Protege) dipper/sources/ClinVar.py -f ClinVarTestSet.xml.gz -o ClinVarTestSet\_‘datestamp’.nt

For while we are still required to redundantly conflate the owl properties in with the data files.

```
python3      ./scripts/add-properties2turtle.py      --input      ./out/ClinVarTestSet_‘datestamp’.nt      --output  
./out/ClinVarTestSet_‘datestamp’.nt --format nt
```

dipper.sources.ClinVar.**allele\_to\_triples**(allele, triples) → None

Process allele info such as dbsnp ids and synonyms :param allele: Allele :param triples: List, Buffer to store the triples :return: None

dipper.sources.ClinVar.**digest\_id**(wordage)

return a deterministic digest of input the ‘b’ is an experiment forcing the first char to be non numeric but valid hex; which is in no way required for RDF but may help when using the identifier in other contexts which do not allow identifiers to begin with a digit

:param wordage the string to hash :returns 20 hex char digest

dipper.sources.ClinVar.**expand\_curie**(this\_curie)

dipper.sources.ClinVar.**is\_literal**(thing)

make inference on type (literal or CURIE)

return: logical

dipper.sources.ClinVar.**make\_biolink\_category\_triple**(subj, cat)

dipper.sources.ClinVar.**make\_spo**(sub, prd, obj, subject\_category=None, object\_category=None)

Decorates the three given strings as a line of ntriples (also writes a triple for subj biolink:category and obj biolink:category)

dipper.sources.ClinVar.**parse**()

Main function for parsing a clinvar XML release and outputting triples

```
dipper.sources.ClinVar.process_measure_set(measure_set,      rcv_acc)      →      dip-  
per.models.ClinVarRecord.Variant
```

Given a MeasureSet, create a Variant object :param measure\_set: XML object :param rcv\_acc: str rcv accession :return: Variant object

```
dipper.sources.ClinVar.record_to_triples(rcv: dipper.models.ClinVarRecord.ClinVarRecord,  
                                         triples: List[T], g2p_map: Dict[KT, VT]) →  
                                         None
```

Given a ClinVarRecord, adds triples to the triples list

### Parameters

- **rcv** – ClinVarRecord
- **triples** – List, Buffer to store the triples
- **g2p\_map** – Gene to phenotype dict

### Returns

None

```
dipper.sources.ClinVar.resolve(label)
composite mapping given f(x) and g(x) here: GLOBALTT & LOCALTT respectivly in order of preference
return g(f(x))|f(x)|g(x) | x TODO consider returning x on fall through

# the decendent resolve(label) function in Source.py # should be used instead and this f(x) removed

: return label's mapping

dipper.sources.ClinVar.csv_link(scv_sig, rcv_trip)
Creates links between SCV based on their pathonicty/significance calls

# GENO:0000840 - GENO:0000840 -> is_equivalent_to SEPIO:0000098 # GENO:0000841 - GENO:0000841
-> is_equivalent_to SEPIO:0000098 # GENO:0000843 - GENO:0000843 -> is_equivalent_to SEPIO:0000098
# GENO:0000844 - GENO:0000844 -> is_equivalent_to SEPIO:0000098 # GENO:0000840 - GENO:0000844
-> contradicts SEPIO:0000101 # GENO:0000841 - GENO:0000844 -> contradicts SEPIO:0000101 #
GENO:0000841 - GENO:0000843 -> contradicts SEPIO:0000101 # GENO:0000840 - GENO:0000841
-> is_consistent_with SEPIO:0000099 # GENO:0000843 - GENO:0000844 -> is_consistent_with SE-
PIO:0000099 # GENO:0000840 - GENO:0000843 -> strongly contradicts SEPIO:0000100

dipper.sources.ClinVar.write_review_status_scores()
Make triples that attach a "star" score to each of ClinVar's review statuses. (Stars are basically a 0-4 rating of the review status.)

Per https://www.ncbi.nlm.nih.gov/clinvar/docs/details/ Table 1. The review status and assignment of stars( with changes made mid-2015) Number of gold stars Description and review statuses

NO STARS: <ReviewStatus> "no assertion criteria provided" <ReviewStatus> "no assertion provided" No submitter provided an interpretation with assertion criteria (no assertion criteria provided), or no interpretation was provided (no assertion provided)

ONE STAR: <ReviewStatus> "criteria provided, single submitter" <ReviewStatus> "criteria provided, conflicting interpretations" One submitter provided an interpretation with assertion criteria (criteria provided, single submitter) or multiple submitters provided assertion criteria but there are conflicting interpretations in which case the independent values are enumerated for clinical significance (criteria provided, conflicting interpretations)

TWO STARS: <ReviewStatus> "criteria provided, multiple submitters, no conflicts" Two or more submitters providing assertion criteria provided the same interpretation (criteria provided, multiple submitters, no conflicts)

THREE STARS: <ReviewStatus> "reviewed by expert panel" reviewed by expert panel

FOUR STARS: <ReviewStatus> "practice guideline" practice guideline A group wishing to be recognized as an expert panel must first apply to ClinGen by completing the form that can be downloaded from our ftp site.

:param None :return: list of triples that attach a "star" score to each of ClinVar's review statuses

dipper.sources.ClinVar.write_spo(sub, prd, obj, triples, subject_category=None, object_category=None)
write triples to a buffer in case we decide to drop them
```

## dipper.sources.Coriell module

```
class dipper.sources.Coriell(graph_type, are_bnodes_skolemized,
                           data_release_version=None)
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.

```
column_labels = ['catalog_id', 'description', 'omim_num', 'sample_type', 'cell_line_availability']  
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 ping 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.yaml file, like: dbauth : {“coriell” : {“user” : “<username>”, “password” : “<password>”, “host” : <host>, “private\_key”=path/to/rsa\_key} }

**Parameters** `is_dl_forced` –

**Returns**

```
files = {'NHGRI': {'columns': ['catalog_id', 'description', 'omim_num', 'sample_type']}
```

```
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_lines = ['ND02380', 'ND02381', 'ND02383', 'ND02384', 'GM17897', 'GM17898', 'GM17899']
```

### dipper.sources.Decipher module

### dipper.sources.EBIGene2Phen module

```
class dipper.sources.EBIGene2Phen(graph_type, are_bnodes_skolemized,  
data_release_version=None)
```

Bases: `dipper.sources.Source`

From EBI: The gene2phenotype dataset (G2P) integrates data on genes, variants and phenotypes for example relating to developmental disorders. It is constructed entirely from published literature, and is primarily an inclusion list to allow targeted filtering of genome-wide data for diagnostic purposes. The dataset was compiled with respect to published genes, and annotated with types of disease-causing gene variants. Each row of the dataset associates a gene with a disease phenotype via an evidence level, inheritance mechanism and mutation consequence. Some genes therefore appear in the database more than once, where different genetic mechanisms result in different phenotypes.

Disclaimer: <https://www.ebi.ac.uk/gene2phenotype/disclaimer> Terms of Use: <https://www.ebi.ac.uk/about/terms-of-use#general> Documentation: <https://www.ebi.ac.uk/gene2phenotype/documentation>

[https://www.clinicalgenome.org/site/assets/files/2757/fitzpatrick\\_ddg2p.pdf](https://www.clinicalgenome.org/site/assets/files/2757/fitzpatrick_ddg2p.pdf)

This script operates on the Developmental Disorders (DDG2P.csv) file In the future we may update to include the cancer gene disease pairs in the CancerG2P.csv file

```
EBI_BASE = 'https://www.ebi.ac.uk/gene2phenotype/downloads/'
```

```
fetch(is_dl_forced: bool = False)
```

Fetch DDG2P.csv.gz and check headers to see if it has been updated

**Parameters** **is\_dl\_forced** – {bool}

**Returns** None

```
files = {'developmental_disorders': {'columns': ['gene_symbol', 'gene_omim_id', 'dis-
```

```
map_files = {'mondo_map': 'https://data.monarchinitiative.org/dipper/cache/unmapped_e-
```

```
parse(limit: Optional[int] = None)
```

Here we parse each row of the gene to phenotype file

We create anonymous variants along with their attributes (allelic requirement, functional consequence) and connect these to genes and diseases

genes are connected to variants via global\_terms[‘has\_affected\_locus’]

variants are connected to attributes via: global\_terms[‘has\_allelic\_requirement’]  
global\_terms[‘has\_functional\_consequence’]

variants are connected to disease based on mappings to the DDD category column, see the translationtable specific to this source for mappings

For cases where there are no disease OMIM id, we either use a disease cache file with mappings to MONDO that has been manually curated

**Parameters** **limit** – {int} number of rows to parse

**Returns** None

## dipper.sources.EOM module

```
class dipper.sources.EOM(graph_type, are_bnodes_skolemized, data_release_version=None)
```

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.yaml 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.

```
GHRRAW = 'https://raw.githubusercontent.com/obophenotype/human-phenotype-ontology'
```

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

connection details for DISCO

```
files = {'map': {'columns': ['morphology_term_id', 'morphology_term_label', 'HP_ID',
```

```
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

resources = {'tables': { 'columns': [ 'morphology_term_id', 'morphology_term_num', 'mo
tables = ['dvp.pr_nlx_157874_1']
```

### dipper.sources.Ensembl module

```
class dipper.sources.Ensembl.Ensembl (graph_type, are_bnodes_skolemized,
                                         data_release_version=None, 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 NCBI Gene and ENSG ids using ENSEMBL's Biomart services.

columns = {'bmq_attributes': ['ensembl_gene_id', 'external_gene_name', 'description'],
fetch (is_dl_forced=True)
    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 mapping from proteins to ensembl_gene(S)? for a species in biomart :param taxid: :return: dict

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.t
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.FlyBase module

```
class dipper.sources.FlyBase.FlyBase (graph_type, are_bnodes_skolemized,
                                         data_release_version=None)
Bases: dipper.sources.PostgreSQLSource.PostgreSQLSource
```

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

Here, we connect to their public database and download preprocessed files

Queries from the relational db 1. allele-phenotype data: ../../sources/sql/fb/allele\_phenotype.sql 2. gene dbxrefs: ../../resources/sql/fb/gene\_xref.sql

Downloads: 1. allele\_human\_disease\_model\_data\_fb\_\*.tsv.gz - models of disease 2. species.ab.gz - species prefix mappings 3. fbal\_to\_fbgn\_fb\*.tsv.gz - allele to gene 4. fbrf\_pmcid\_pmcid\_doi\_fb\_\*.tsv.gz - flybase ref to pmid

We connect using the [Direct Chado Access]([http://gmod.org/wiki/Public\\_Chado\\_Databases#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`.

Note that this script underwent a major revision after commit bd5f555 in which genotypes, stocks, and environments were removed

```
CURREL = 'releases/current/precomputed_files'
FLYFTP = 'ftp.firebaseio.net'

fetch(is_dl_forced=False)
Fetch flat files and sql queries

    Parameters is_dl_forced – force download

    Returns None

files = {'allele_gene': {'columns': ['AlleleID', 'AlleleSymbol', 'GeneID', 'GeneSymbol']}
parse(limit=None)
Parse flybase files and add to graph

    Parameters limit – number of rows to process

    Returns None

queries = {'allele_phenotype': {'columns': ['allele_id', 'pheno_desc', 'pheno_type',

```

## dipper.sources.GWASCatalog module

```
class dipper.sources.GWASCatalog(graph_type,      are_bnodes_skolemized,
                                  data_release_version=None)
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](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': {'columns': ['DATE ADDED TO CATALOG', 'PUBMEDID', 'FIRST AUTHOR']}
```

```
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

### dipper.sources.GeneOntology module

```
class dipper.sources.GeneOntology.GeneOntology(graph_type, are_bnodes_skolemized,
                                                data_release_version=None,
                                                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

**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': {'columns': ['DB', 'DB\_Object\_ID', 'DB\_Object\_Symbol', 'Qualifier']}

**gaf\_columns** = ['DB', 'DB\_Object\_ID', 'DB\_Object\_Symbol', 'Qualifier', 'GO\_ID', 'DB\_Ref']

**getTestSuite**()

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

**get\_uniprot\_entrez\_id\_map**()

**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**(*gaffile, limit, id\_map=None*)

**wont\_prefix** = ['zgc', 'wu', 'si', 'im', 'BcDNA', 'sb', 'anon-EST', 'EG', 'id', 'zmp',

### dipper.sources.GeneReviews module

### dipper.sources.HGNC module

### dipper.sources.HPOAnnotations module

```
class dipper.sources.HPOAnnotations.HPOAnnotations(graph_type,
                                                    are_bnodes_skolemized,
                                                    data_release_version=None)
```

Bases: *dipper.sources.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 resources/test\_ids.yaml 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()

The (several thousands) common-disease files from the repo tarball are added to the files object. try adding the ‘common-disease-mondo’ files as well?

**fetch** (is *dl forced*=*False*)

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

```
files = {'doi':
```

For more information about the study, please contact Dr. John Smith at (555) 123-4567 or via email at [john.smith@researchinstitute.org](mailto:john.smith@researchinstitute.org).

## An abstract method `th`

**common\_files()**  
Each file has its own initialization function. Many of them will be initialized in the main() function, but some may be initialized later.

a tarball

## Returns

**se** (*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

**cess\_all\_common\_disease\_files** (*limit=None*)  
Loop through all of the files that we previously fetched from git, creating the disease-phenotype associa-

tion. :param limit: :return:

**cess\_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

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

Restaurante

```
small_files = {'columns': ['Disease_ID', 'Disease_Name', 'Gene_ID', 'Gene_Name', 'Gene_Symbol']}
```

#### dipper sources | MPC module

```
class dipper.sources.IMPc.IMPc(graph_type,  
                                data_release_version=None)  
Bases: dipper_sources.Source, Source
```

From the [IMPC](<https://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 = {'checksum': {'file': 'genotype-phenotype-assertions-ALL.csv.tgz.md5', 'url':**

### **getTestSuite ()**

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

### **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**

### **parse\_checksum\_file (file)**

:param file :return dict

## dipper.sources.KEGG module

## dipper.sources.MGI module

```
class dipper.sources.MGI.MGI(graph_type, are_bnodes_skolemized, data_release_version=None)
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.yaml file, like: dbauth : {‘mgi’ : {‘user’ : ‘<username>’, ‘password’ : ‘<password>’}} You can request access by contacting [mgi-help@jax.org](mailto: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**

**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 rows in 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_map':  [{‘query’: ‘.../..../resources/sql/mgi/mgi_dbinfo.sql’, ‘outf...’}], 
tables = {‘all_allele_mutation_view’: {‘columns’: [‘_allele_key’, ‘mutation’]}}, 
unknown_taxa = [‘Not Applicable’, ‘Not Specified’]}
```

## dipper.sources.MGISlim module

```
class dipper.sources.MGISlim.MGISlim(graph_type, are_bnodes_skolemized, data_release_version=None)
Bases: dipper.sources.Source.Source
```

slim mgi model only containing Gene to phenotype associations Uses mousemine: <http://www.mousemine.org/mousemine/begin.do> python lib api <http://intermine.org/intermine-ws-python/>

**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

## dipper.sources.MMRRC module

```
class dipper.sources.MMRRC(graph_type, are_bnodes_skolemized, data_release_version=None)
Bases: dipper.sources.Source
```

Here we process the Mutant Mouse Resource and Research Center (<https://www.mmrcc.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 MMRCC 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': {'columns': ['STRAIN/STOCK_ID', 'STRAIN/STOCK_DESIGNATION', 'STRAN'}}
```

### **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-U'
```

## dipper.sources.MPD module

```
class dipper.sources.MPD(graph_type, are_bnodes_skolemized, data_release_version=None)
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 +/- 2 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, localtt)

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': {'columns': ['measnum', 'mpdsector', 'projsym', 'varname']}
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'

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,
                               data_release_version=None)
Bases: dipper.sources.Source.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

files = {'omia_d2p': {'columns': ['Disease ID', 'Species ID', 'Breed Name', 'Variant
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,
                                 data_release_version=None, 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. (note: in truth we create blank nodes and then pretend they are something else. TEC)

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 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': {'build_num': 'mm10', 'file': '10090cytoBand.txt.gz', 'genome_lab': 'mm10'}}
```

**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

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

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,
                           data_release_version=None)
Bases: dipper.sources.Source.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,
                           data_release_version=None)
Bases: dipper.sources.Source.Source

Drugs and Compounds stored in the BioThings database
MY_DRUG_API = 'http://c.biotothings.io/v1/query'
check_if_remote_is_newer(localfile)
    Need to figure out how biotothings 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

### dipper.sources.OMIA module

### dipper.sources.OMIM module

### dipper.sources.OMIMSource module

### dipper.sources.Orphanet module

```
class dipper.sources.Orphanet(graph_type, are_bnodes_skolemized,  
                                data_release_version=None)
```

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.

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

**Parameters** **is\_dl\_forced** –

**Returns**

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

```
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.Panther module

```
class dipper.sources.Panther(graph_type, are_bnodes_skolemized,  
                                data_release_version=None, tax_ids=None)
```

Bases: *dipper.sources.Source*

The pairwise orthology calls from Panther DB: <http://pantherdb.org/> encompass 22 species, from the RefGenome 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 resources/test\_id.yaml.

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 = {'Orthologs_HCOP': {'columns': ['Gene', 'Ortholog', 'Type of ortholog', 'Comm
```

```
getTestSuite()
```

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

```
panther_format = ['Gene', 'Ortholog', 'Type of ortholog', 'Common ancestor for the ort
```

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

**Returns** None

## dipper.sources.PostgreSQLSource module

```
class dipper.sources.PostgreSQLSource.PostgreSQLSource(graph_type,
are_bnodes_skolemized,
data_release_version=None,
name=None,                  in-
gest_title=None,           in-
gest_url=None,            in-
gest_logo=None,           in-
gest_description=None,
license_url=None,
data_rights=None,
file_handle=None)
```

Bases: `dipper.sources.Source.Source`

Class for interfacing with remote Postgres databases

```
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_pgdb(tables, cxn, limit=None)
```

**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)
```

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:

```
files = {}
```

```
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

### dipper.sources.RGD module

```
class dipper.sources.RGD (graph_type, are_bnodes_skolemized, data_release_version=None)  
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 via the rat\_genome\_database ftp site  
Args:

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

**Returns:** :return None

```
files = {'rat_gene2mammalian_phenotype': {'columns': ['DB', 'DB Object ID', 'DB Objec'}}
```

```
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,  
                                         data_release_version=None)
```

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': {'columns': ['component', 'pathway_id', 'pathway_iri', 'pa
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, data_release_version=None)
    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 yeast_genome_database using the
        yeast_genome_download site.

    Args:
        param is_dl_forced (bool) Force download

    Returns: :return None

    files = {'sgd_phenotype': {'columns': ['Feature Name', 'Feature Type', 'Gene Name',
    static make_apo_map()

    make_association (record)
        construct the association :param record: :return: modeled association of genotype to mammalian??? phe-
        notype

    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='rdf_graph', are_bnodes_skized=False,
                            data_release_version=None, name=None, in-
                            gest_title=None, ingest_url=None, ingest_logo=None,
                            ingest_description=None, license_url=None,
                            data_rights=None, file_handle=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.

Also provides a means to marshal metadata in a consistent fashion

Houses the global translation table (from ontology label to ontology term) so it may as well be used everywhere.

```
ARGV = {}

DIPPERCACHE = 'https://archive.monarchinitiative.org/DipperCache'

static check_fileheader(expected, received, src_key=None)
    Compare file headers received versus file headers expected if the expected headers are a subset (proper or not) of received headers report success (warn if proper subset)

        param: expected list param: received list

        return: truthyness

check_if_remote_is_newer(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

command_args()
    To make arbitrary variables from dipper-etl.py's calling environment available when working in source ingestors in a hopefully universal way

    Does not appear to be populated till after an ingest's _init_() finishes.

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

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(remoteurl, 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 remoteurl: URL of remote file to fetch :param localfile: pathname of file to save locally

        Returns bool

static file_len(fname)
files = {}

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

static get_file_md5(directory, filename, blocksize=1048576)
get_files(is_dl_forced, files=None, delay=0)
    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

static 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(wordage)
prepend 'b' to avoid leading with digit truncate to a 20 char sized word with a leading 'b' return truncated
sha1 hash of string.

by the birthday paradox; expect 50% chance of collision after 69 billion invocations however these are
only hoped to be unique within a single file

Consider reducing to 17 hex chars to fit in a 64 bit word 16 discounting a leading constant gives a 50%
chance of collision at about 4.3b billion unique input strings (currently _many_ orders of magnitude below
that)

Parameters long_string – str string to be hashed

Returns str hash of id

load_local_translationtable(name)

Load “ingest specific” translation from whatever they called something to the ontology label we
need to map it to. To facilitate seeing more ontology labels in dipper ingests a reverse mapping
from ontology labels to external strings is also generated and available as a dict localcid

‘—
# %s.yaml “”: “” # example’

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(yamlfile)

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 dump of a mysql relational database. The
“elem” is akin to a mysql table, with it’s name of `table_name` . It will process each `row` given the
`processing_function` supplied. :param elem: The element data :param table_name: The name
of the table to process :param processing_function: The row processing function :param 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**

**resolve** (`word`, `mandatory=True`, `default=None`)

composite mapping given  $f(x)$  and  $g(x)$  here: `localtt` & `globaltt` respectively return  $g(f(x))|g(x)|f(x)|x$  in order of preference returns  $x|default$  on fall through if finding a mapping is not mandatory (by default finding is mandatory).

This may be specialized further from any mapping to a global mapping only; if need be.

**Parameters**

- **word** – the string to find as a key in translation tables
- **mandatory** – boolean to cause failure when no key exists
- **default** – string to return if nothing is found (& not mandatory)

**:return** value from global translation table, or value from local translation table, or the query key if finding a value is not mandatory (in this order)

**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** ()

pointless convieniance

**write** (`fmt='turtle'`, `stream=None`, `write_metadata_in_main_graph=True`)

**This convenience method will write out all of the graphs** associated with the source.

Right now these are hardcoded to be a single main “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,  
                                         data_release_version=None, 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](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>

**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(graph_type, are_bnodes_skolemized,
                                data_release_version=None, 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

```
rdf: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 its 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'
cytobandideo_columns = ['chrom', 'chromStart', 'chromEnd', 'name', 'giestain']
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': {'assembly': ('UCSC:mm10', 'UCSC:mm9'), 'build_num': 'mm10', 'col...
```

`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.UDP(graph_type, are_bnodes_skolemized, data_release_version=None)  
    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,
                                         data_release_version=None)
Bases: dipper.sources.Source.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': {'columns': ['DB', 'DB Object ID', 'DB Object Symbol', 'Qua'}
```

**static 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'
update_wsnr_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'
wbre1 = 'ftp://ftp.wormbase.org/pub/wormbase/releases'
```

### dipper.sources.Xenbase module

```
class dipper.sources.Xenbase.Xenbase(graph_type, are_bnodes_skolemized,
                                      data_release_version=None)
Bases: dipper.sources.Source.Source
Xenbase is a web-accessible resource that integrates all the diverse biological, genomic, genotype and phenotype data available from Xenopus research.

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_assertions': {'columns': ['SUBJECT', 'SUBJECT_LABEL', 'SUBJECT_TAXON']}
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.ZFIN module

```
class dipper.sources.ZFIN.ZFIN(graph_type, are_bnodes_skolemized,
                                 data_release_version=None)
Bases: dipper.sources.Source.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).

Genotype labels are output as ZFIN genotype name + “[background]”. We also process the genotype components to build labels in a monarch-style, and these are added as synonyms. The monarch-style genotype label includes: \* 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>/+

see: resources/zfin/README for column extraction from downloads page

### **fetch (is\_dl\_forced=False)**

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

```
fhandle = <_io.TextIOWrapper name='/home/docs/checkouts/readthedocs.org/user_builds/dipper/
```

```
files = {'backgrounds': { 'columns': [ 'Genotype ID', 'Genotype Name', 'Background', ''] }}
```

### **static get\_orthology\_evidence\_code (abbrev)**

move to localtt & globltt

### **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)**

```
test_ids = { 'allele': [ 'ZDB-ALT-010426-4', 'ZDB-ALT-010427-8', 'ZDB-ALT-011017-8', 'ZDB-ALT-011018-8' ] }
```

## dipper.sources.ZFINSlim module

```
class dipper.sources.ZFINSlim(graph_type, are_bnodes_skolemized, data_release_version=None)
Bases: dipper.sources.Source
```

zfin mgi model only containing Gene to phenotype associations Using the file here: [https://zfin.org/downloads/phenoGeneCleanData\\_fish.txt](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': {'columns': ['ID', 'Gene Symbol', 'Gene ID', 'Affected Structure'],
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(pxf)
```

##### 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_hgnc_id_from_symbol(gene_symbol)
        Get HGNC curie from symbol using monarch and mygene services :param gene_symbol: :return:

    static get_homologene_by_gene_num(gene_num)
    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_id_in_mondo(curie, mondo_min)
```

**Parameters**

- **curie** – curie formatted ID
- **mondo\_min** – a json decoded mondo-minimal.json file, eg <https://github.com/monarch-initiative/mondo/releases/> download/2019-04-06/mondo-minimal.json

**Returns** boolean, true if ID is in mondo and false otherwise

```
static remove_control_characters(string)
```

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(curie_map)
```

Bases: object

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

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

```
static compare_graph_predicates(graph1, graph2)
```

From rdf graphs, count predicates in each and return a list of : param graph1 graph, hopefully RDFLib-like : param graph2 graph, ditto : return dict with count of predicates in each graph: : e.g.: : { : “has\_a\_property”: { : “graph1”: 1234, : “graph2”: 1023}, : “has\_another\_property”: { : “graph1”: 94, : “graph2”: 51} : }

```
static count_predicates(graph)
```

From rdf graphs, count predicates in each and return a list of : param graph : return dict with count of predicates in each graph: : e.g.: : { : “has\_a\_property”: 1234, : “has\_another\_property”: 482 : }

```
static digest_id(wordage)
```

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

```
static write(graph, fileformat=None, filename=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

```
static remove_ontology_axioms(graph)
```

Given an rdflib graph, remove any triples connected to an ontology node: {} a owl:Ontology :param graph: RDFGraph :return: None

```
static test_graph_equality(turtlish, graph)
```

### Parameters

- **turtlish** – file path or 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.rdf2dot module

A fork of rdflib rdf2dot utility, see [https://rdflib.readthedocs.io/en/stable/\\_modules/rdflib/tools/rdf2dot.html#rdf2dot](https://rdflib.readthedocs.io/en/stable/_modules/rdflib/tools/rdf2dot.html#rdf2dot)

We apply the formatliteral function to labels, but otherwise the code is the same

This is necessary for variants with HGVS primary labels that contain characters that need to be url encoded (<, >)

Also replaces cgi.escape with html.escape

TO DO make a PR

`dipper.utils.rdf2dot.rdf2dot(g, stream, graph_opts=None)`

Convert the RDF graph to DOT writes the dot output to the stream

## dipper.utils.romanplus module

Convert to and from Roman numerals This program is part of “Dive Into Python”, a free Python tutorial for experienced programmers. Visit <http://diveintopython.org/> for the latest version.

This program is free software; you can redistribute it and/or modify it under the terms of the Python 2.1.1 license, available at <http://www.python.org/2.1.1/license.html>

Note: This has been modified to add optional characters after the initial roman numbers by nlw.

`dipper.utils.romanplus.fromRoman(strng)`  
convert Roman numeral to integer

`dipper.utils.romanplus.toRoman(num)`  
convert integer to Roman numeral

## Submodules

### dipper.config module

`dipper.config.conf = {'keys': {'omim': ''}}`  
Load the configuration file ‘conf.yaml’, if it exists. it isn’t always required, but may be for some sources.  
conf.yaml 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, 64  
dipper.curie\_map, 64  
dipper.graph, 13  
dipper.graph.Graph, 13  
dipper.graph.RDFGraph, 13  
dipper.graph.StreamedGraph, 14  
dipper.models, 15  
dipper.models.assoc, 15  
dipper.models.assoc.Association, 15  
dipper.models.assoc.Chem2DiseaseAssoc, 16  
dipper.models.assoc.D2PAssoc, 16  
dipper.models.assoc.G2PAssoc, 17  
dipper.models.assoc.InteractionAssoc, 18  
dipper.models.assoc.OrthologyAssoc, 18  
dipper.models.BiolinkVocabulary, 18  
dipper.models.ClinVarRecord, 18  
dipper.models.Dataset, 19  
dipper.models.Environment, 23  
dipper.models.Evidence, 24  
dipper.models.Family, 25  
dipper.models.GenomicFeature, 25  
dipper.models.Genotype, 26  
dipper.models.Model, 29  
dipper.models.Pathway, 31  
dipper.models.Provenance, 32  
dipper.models.Reference, 32  
dipper.sources, 33  
dipper.sources.AnimalQTlDb, 33  
dipper.sources.Bgee, 34  
dipper.sources.BioGrid, 34  
dipper.sources.ClinVar, 35  
dipper.sources.Coriell, 37  
dipper.sources.CTD, 35  
dipper.sources.EBIGene2Phen, 38  
dipper.sources.Ensembl, 40

dipper.sources.EOM, 39  
dipper.sources.FlyBase, 40  
dipper.sources.GeneOntology, 42  
dipper.sources.GWASCatalog, 41  
dipper.sources.HPOAnnotations, 42  
dipper.sources.IMPC, 43  
dipper.sources.MGI, 45  
dipper.sources.MGISlim, 45  
dipper.sources.MMRRC, 46  
dipper.sources.Monarch, 47  
dipper.sources.Monochrom, 47  
dipper.sources.MPD, 46  
dipper.sources.MyChem, 49  
dipper.sources.MyDrug, 49  
dipper.sources.Orphanet, 50  
dipper.sources.Panther, 50  
dipper.sources.PostgreSQLSource, 51  
dipper.sources.Reactome, 52  
dipper.sources.RGD, 52  
dipper.sources.SGD, 53  
dipper.sources.Source, 53  
dipper.sources.StringDB, 56  
dipper.sources.UCSCBands, 57  
dipper.sources.UDP, 58  
dipper.sources.WormBase, 59  
dipper.sources.Xenbase, 60  
dipper.sources.ZFIN, 60  
dipper.sources.ZFINSlim, 61  
dipper.utils, 62  
dipper.utils.CurieUtil, 62  
dipper.utils.DipperUtil, 62  
dipper.utils.GraphUtils, 63  
dipper.utils.rdf2dot, 64  
dipper.utils.romanplus, 64  
dipper.utils.TestTools, 63



---

## Index

---

### A

add_agent_to_graph()	(dip- per.models.Provenance.Provenance method),	15
add_assay_to_graph()	(dip- per.models.Provenance.Provenance method),	63
add_assertion()	(dip- per.models.Provenance.Provenance method),	63
add_association_to_graph()	(dip- per.models.assoc.Association.Assoc method),	15
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),	43
add_data_individual()	(dip- per.models.Evidence.Evidence method),	24
add_date()	(dipper.models.assoc.Association.Assoc method),	15
add_date_created()	(dip- per.models.Provenance.Provenance method),	32
add_evidence()	(dip- per.models.assoc.Association.Assoc method),	24
add_evidence()	(dipper.models.Evidence.Evidence method),	24
add_gene_family_to_graph()	(dip- per.models.assoc.OrthologyAssoc.OrthologyAssoc method),	25
add_predicate_object()	(dip-	
	per.models.assoc.Association.Assoc method),	15
	add_property_axioms()	(dip- per.utils.GraphUtils.GraphUtils static method),
	add_property_to_graph()	(dip- per.utils.GraphUtils.GraphUtils static method),
	add_provenance()	(dip- per.models.assoc.Association.Assoc method),
	add_relation()	(dipper.sources.MyChem.MyChem static method),
	add_source()	(dip- per.models.assoc.Association.Assoc method),
	add_source()	(dipper.models.Evidence.Evidence method),
	add_study_measure()	(dip- per.models.Provenance.Provenance method),
	add_study_parts()	(dip- per.models.Provenance.Provenance method),
	add_study_to_measurements()	(dip- per.models.Provenance.Provenance method),
	add_supporting_data()	(dip- per.models.Evidence.Evidence method),
	add_supporting_evidence()	(dip- per.models.Evidence.Evidence method),
	add_supporting_publication()	(dip- per.models.Evidence.Evidence method),
	addAffectedLocus()	(dip- per.models.Genotype.Genotype method),
	addAllele()	(dipper.models.Genotype.Genotype

<i>method), 27</i>		<i>25</i>
<code>addAlleleOfGene ()</code> <i>per.models.Genotype.Genotype</i> <i>27</i>	<i>(dip-</i> <i>method),</i>	<code>addFeatureProperty ()</code> <i>per.models.GenomicFeature.Feature</i> <i>method),</i> <i>25</i>
<code>addAuthor ()</code> <i>(dipper.models.Reference.Reference</i> <i>method), 32</i>		<code>addFeatureStartLocation ()</code> <i>(dip-</i> <i>per.models.GenomicFeature.Feature</i> <i>method),</i> <i>25</i>
<code>addBlankNodeAnnotation ()</code> <i>per.models.Model.Model</i> <i>method), 29</i>	<i>(dip-</i> <i>method),</i>	<code>addFeatureToGraph ()</code> <i>(dip-</i> <i>per.models.GenomicFeature.Feature</i> <i>method),</i> <i>25</i>
<code>addChromosome ()</code> <i>per.models.Genotype.Genotype</i> <i>27</i>	<i>(dip-</i> <i>method),</i>	<code>addGene ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method), 27</i>
<code>addChromosomeClass ()</code> <i>per.models.Genotype.Genotype</i> <i>27</i>	<i>(dip-</i> <i>method),</i>	<code>addGeneProduct ()</code> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>27</i>
<code>addChromosomeInstance ()</code> <i>per.models.Genotype.Genotype</i> <i>27</i>	<i>(dip-</i> <i>method),</i>	<code>addGeneTargetingReagent ()</code> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addClassToGraph ()</code> <i>(dipper.models.Model.Model</i> <i>method), 29</i>		<code>addGeneTargetingReagentToGenotype ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addComment ()</code> <i>(dipper.models.Model.Model</i> <i>method),</i> <i>30</i>		<code>addGeneToPathway ()</code> <i>(dip-</i> <i>per.models.Pathway.Pathway</i> <i>method), 31</i>
<code>addComponentAttributes ()</code> <i>per.models.Environment.Environment</i> <i>method),</i> <i>23</i>	<i>(dip-</i> <i>per.models.Environment.Environment</i> <i>method),</i>	<code>addGenome ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addComponentToEnvironment ()</code> <i>per.models.Environment.Environment</i> <i>method),</i> <i>23</i>	<i>(dip-</i> <i>per.models.Environment.Environment</i> <i>method),</i>	<code>addGenomicBackground ()</code> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addComponentToPathway ()</code> <i>per.models.Pathway.Pathway</i> <i>method), 31</i>	<i>(dip-</i> <i>per.models.Pathway.Pathway</i> <i>method),</i>	<code>addGenomicBackgroundToGenotype ()</code> <i>(dip-</i> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addConstruct ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method), 27</i>		<code>addGenotype ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addDefinition ()</code> <i>(dipper.models.Model.Model</i> <i>method), 30</i>	<i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i>	<code>addIndividualToGraph ()</code> <i>(dip-</i> <i>per.models.Model.Model</i> <i>method), 30</i>
<code>addDepiction ()</code> <i>(dipper.models.Model.Model</i> <i>method), 30</i>		<code>addLabel ()</code> <i>(dipper.models.Model.Model</i> <i>method),</i> <i>30</i>
<code>addDeprecatedClass ()</code> <i>per.models.Model.Model</i> <i>method), 30</i>	<i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i>	<code>addMember ()</code> <i>(dipper.models.Family.Family</i> <i>method),</i> <i>25</i>
<code>addDeprecatedIndividual ()</code> <i>per.models.Model.Model</i> <i>method), 30</i>	<i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i>	<code>addMemberOf ()</code> <i>(dipper.models.Family.Family</i> <i>method),</i> <i>25</i>
<code>addDerivesFrom ()</code> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>27</i>	<i>(dip-</i> <i>per.models.Genotype.Genotype</i> <i>method),</i>	<code>addMemberOfPopulation ()</code> <i>per.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>
<code>addDescription ()</code> <i>(dipper.models.Model.Model</i> <i>method), 30</i>		<code>addOntologyDeclaration ()</code> <i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i> <i>30</i>
<code>addEnvironment ()</code> <i>per.models.Environment.Environment</i> <i>method),</i> <i>23</i>		<code>addOWLPropertyClassRestriction ()</code> <i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i> <i>30</i>
<code>addEnvironmentalCondition ()</code> <i>per.models.Environment.Environment</i> <i>method),</i> <i>23</i>		<code>addOWLVersionInfo ()</code> <i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i> <i>30</i>
<code>addEquivalentClass ()</code> <i>per.models.Model.Model</i> <i>method), 30</i>	<i>(dip-</i> <i>per.models.Model.Model</i> <i>method),</i>	<code>addPage ()</code> <i>(dipper.models.Reference.Reference</i> <i>method),</i> <i>32</i>
<code>addFeatureEndLocation ()</code> <i>per.models.GenomicFeature.Feature</i> <i>method),</i>	<i>(dip-</i> <i>per.models.GenomicFeature.Feature</i> <i>method),</i>	<code>addParts ()</code> <i>(dipper.models.Genotype.Genotype</i> <i>method),</i> <i>28</i>

addPartsToVSLC ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>addTitle ()</i> ( <i>dipper.models.Reference.Reference method</i> ), 32
<i>per.models.Genotype.Genotype</i>		
28		
addPathway ()	( <i>dipper.models.Pathway.Pathway method</i> ), 31	<i>addTriple ()</i> ( <i>dipper.graph.Graph.Graph method</i> ), 13
addPerson ()	( <i>dipper.models.Model.Model method</i> ),	<i>addTriple ()</i> ( <i>dipper.graph.RDFGraph.RDFGraph method</i> ), 14
30		
addPolypeptide ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>addTriple ()</i> ( <i>dipper.models.Model.Model method</i> ),
<i>per.models.Genotype.Genotype</i>		31
28		
addPositionToGraph ()	( <i>dipper.models.GenomicFeature.Feature method</i> ),	<i>addType ()</i> ( <i>dipper.models.Model.Model method</i> ), 31
<i>per.models.GenomicFeature.Feature</i>		
26		
addReagentTargetedGene ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>addVSLCtoParent ()</i> ( <i>dipper.models.Genotype.Genotype method</i> ),
<i>per.models.Genotype.Genotype</i>		29
28		
addReferenceGenome ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>addXref ()</i> ( <i>dipper.models.Model.Model method</i> ), 31
<i>per.models.Genotype.Genotype</i>		
29		
addRefToGraph ()	( <i>dipper.models.Reference.Reference method</i> ),	<i>Allele</i> ( <i>class in dipper.models.ClinVarRecord</i> ), 18
<i>per.models.Reference.Reference</i>		
32		
addRegionPositionToGraph ()	( <i>dipper.models.GenomicFeature.Feature method</i> ),	<i>allele_to_triples ()</i> ( <i>in module dipper.sources.ClinVar</i> ), 36
<i>per.models.GenomicFeature.Feature</i>		
26		
addSameIndividual ()	( <i>dipper.models.Model.Model method</i> ),	<i>AnimalQTLdb</i> ( <i>class in dipper.sources.AnimalQTLdb</i> ),
<i>per.models.Model.Model</i>		33
31		
addSequenceAlteration ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>ARGV</i> ( <i>dipper.sources.Source.Source attribute</i> ), 53
<i>per.models.Genotype.Genotype</i>		
29		
addSequenceAlterationToVariantLocus ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>Assoc</i> ( <i>class in dipper.models.assoc.Association</i> ), 15
<i>per.models.Genotype.Genotype</i>		
29		
addSequenceDerivesFrom ()	( <i>dipper.models.Genotype.Genotype method</i> ),	
<i>per.models.Genotype.Genotype</i>		
29		
addSubClass ()	( <i>dipper.models.Model.Model method</i> ),	<b>B</b>
31		
addSubsequenceOfFeature ()	( <i>dipper.models.GenomicFeature.Feature method</i> ),	<i>Bgee</i> ( <i>class in dipper.sources.Bgee</i> ), 34
<i>per.models.GenomicFeature.Feature</i>		
26		
addSynonym ()	( <i>dipper.models.Model.Model method</i> ),	<i>bind_all_namespaces ()</i> ( <i>dipper.graph.RDFGraph.RDFGraph method</i> ),
31		14
addTargetedGeneComplement ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>BioGrid</i> ( <i>class in dipper.sources.BioGrid</i> ), 34
<i>per.models.Genotype.Genotype</i>		
29		
addTargetedGeneSubregion ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>biogrid_ids</i> ( <i>dipper.sources.BioGrid.BioGrid attribute</i> ), 34
<i>per.models.Genotype.Genotype</i>		
29		
addTaxon ()	( <i>dipper.models.Genotype.Genotype method</i> ),	<i>BioLinkVocabulary</i> ( <i>class in dipper.models.BiolinkVocabulary</i> ), 18
29		
addTaxonToFeature ()	( <i>dipper.models.GenomicFeature.Feature method</i> ),	<i>bl_file_with_path</i> ( <i>dipper.models.BiolinkVocabulary.BioLinkVocabulary attribute</i> ), 18
<i>per.models.GenomicFeature.Feature</i>		
26		
		<i>bl_vocab</i> ( <i>dipper.models.BiolinkVocabulary.BioLinkVocabulary attribute</i> ), 18
		<i>build_measurement_description ()</i> ( <i>dipper.sources.MPD.MPD static method</i> ), 47
		<b>C</b>
		<i>check_fileheader ()</i> ( <i>dipper.sources.Source.Source static method</i> ),
		54
		<i>check_if_remote_is_newer ()</i> ( <i>dipper.sources.Bgee.Bgee method</i> ), 34
		<i>check_if_remote_is_newer ()</i> ( <i>dipper.sources.MyDrug.MyDrug method</i> ), 49
		<i>check_if_remote_is_newer ()</i> ( <i>dipper.sources.Source.Source method</i> ), 54
		<i>check_uniprot ()</i> ( <i>dipper.sources.MyChem.MyChem static method</i> ),
		49

Chem2DiseaseAssoc (class in dipper.models.assoc.Chem2DiseaseAssoc), 16  
chunks() (dipper.sources.MyChem.MyChem static method), 49  
ClinVarRecord (class in dipper.models.ClinVarRecord), 19  
column\_labels (dipper.sources.Coriell.Coriell attribute), 38  
columns (dipper.sources.Ensembl.Ensembl attribute), 40  
command\_args() (dipper.sources.Source.Source method), 54  
compare\_checksums() (dipper.sources.IMPC.IMPC method), 44  
compare\_graph\_predicates() (dipper.utils.GraphUtils.GraphUtils static method), 63  
compare\_local\_remote\_bytes() (dipper.sources.Source.Source method), 54  
Condition (class in dipper.models.ClinVarRecord), 19  
conf (in module dipper.config), 64  
Coriell (class in dipper.sources.Coriell), 37  
count\_predicates() (dipper.utils.GraphUtils.GraphUtils static method), 63  
CTD (class in dipper.sources.CTD), 35  
curie\_map (dipper.graph.RDFGraph.RDFGraph attribute), 14  
curie\_map (dipper.graph.StreamedGraph.StreamedGraph attribute), 14  
curie\_regexp (dipper.graph.Graph.Graph attribute), 13  
curie\_util (dipper.graph.RDFGraph.RDFGraph attribute), 14  
curie\_util (dipper.graph.StreamedGraph.StreamedGraph attribute), 14  
CurieUtil (class in dipper.utils.CurieUtil), 62  
CURREL (dipper.sources.FlyBase.FlyBase attribute), 41  
cytobandideo\_columns (dipper.sources.UCSCBands.UCSCBands attribute), 57  
D  
D2PAssoc (class in dipper.models.assoc.D2PAssoc), 16  
Dataset (class in dipper.models.Dataset), 19  
default\_species (dipper.sources.Bgee.Bgee attribute), 34  
digest\_id() (dipper.utils.GraphUtils.GraphUtils static method), 63  
digest\_id() (in module dipper.sources.ClinVar), 36  
dipper (module), 13  
dipper.config (module), 64  
dipper.curie\_map (module), 64  
dipper.graph (module), 13  
dipper.graph.Graph (module), 13  
dipper.graph.RDFGraph (module), 13  
dipper.graph.StreamedGraph (module), 14  
dipper.models (module), 15  
dipper.models.assoc (module), 15  
dipper.models.assoc.Association (module), 15  
dipper.models.assoc.Chem2DiseaseAssoc (module), 16  
dipper.models.assoc.D2PAssoc (module), 16  
dipper.models.assoc.G2PAssoc (module), 17  
dipper.models.assoc.InteractionAssoc (module), 18  
dipper.models.assoc.OrthologyAssoc (module), 18  
dipper.models.BiolinkVocabulary (module), 18  
dipper.models.ClinVarRecord (module), 18  
dipper.models.Dataset (module), 19  
dipper.models.Environment (module), 23  
dipper.models.Evidence (module), 24  
dipper.models.Family (module), 25  
dipper.models.GenomicFeature (module), 25  
dipper.models.Genotype (module), 26  
dipper.models.Model (module), 29  
dipper.models.Pathway (module), 31  
dipper.models.Provenance (module), 32  
dipper.models.Reference (module), 32  
dipper.sources (module), 33  
dipper.sources.AnimalQTLdb (module), 33  
dipper.sources.Bgee (module), 34  
dipper.sources.BioGrid (module), 34  
dipper.sources.ClinVar (module), 35  
dipper.sources.Coriell (module), 37  
dipper.sources.CTD (module), 35  
dipper.sources.EBIGene2Phen (module), 38  
dipper.sources.Ensembl (module), 40  
dipper.sources.EOM (module), 39  
dipper.sources.FlyBase (module), 40  
dipper.sources.GeneOntology (module), 42  
dipper.sources.GWASCatalog (module), 41  
dipper.sources.HPOAnnotations (module), 42  
dipper.sources.IMPC (module), 43  
dipper.sources.MGI (module), 45  
dipper.sources.MGISlim (module), 45  
dipper.sources.MMRRC (module), 46  
dipper.sources.Monarch (module), 47  
dipper.sources.Monochrom (module), 47  
dipper.sources.MPD (module), 46  
dipper.sources.MyChem (module), 49  
dipper.sources.MyDrug (module), 49  
dipper.sources.Orphanet (module), 50  
dipper.sources.Panther (module), 50

dipper.sources.PostgreSQLSource (*module*), 51  
dipper.sources.Reactome (*module*), 52  
dipper.sources.RGD (*module*), 52  
dipper.sources.SGD (*module*), 53  
dipper.sources.Source (*module*), 53  
dipper.sources.StringDB (*module*), 56  
dipper.sources.UCSCBands (*module*), 57  
dipper.sources.UDP (*module*), 58  
dipper.sources.WormBase (*module*), 59  
dipper.sources.Xenbase (*module*), 60  
dipper.sources.ZFIN (*module*), 60  
dipper.sources.ZFINSlim (*module*), 61  
dipper.utils (*module*), 62  
dipper.utils.CurieUtil (*module*), 62  
dipper.utils.DipperUtil (*module*), 62  
dipper.utils.GraphUtils (*module*), 63  
dipper.utils.rdf2dot (*module*), 64  
dipper.utils.romanplus (*module*), 64  
dipper.utils.TestToolss (*module*), 63  
DIPPERCACHE (*dipper.sources.Source.Source tribute*), 54  
DipperUtil (*class in dipper.utils.DipperUtil*), 62

**E**

EBI\_BASE (*dipper.sources.EBIGene2Phen.EBIGene2Phen attribute*), 39  
EBIGene2Phen (*class in dipper.sources.EBIGene2Phen*), 38  
Ensembl (*class in dipper.sources.Ensembl*), 40  
Environment (*class in dipper.models.Environment*), 23  
EOM (*class in dipper.sources.EOM*), 39  
Evidence (*class in dipper.models.Evidence*), 24  
execute\_query () (*dipper.sources.MyChem.MyChem static method*), 49  
expand\_curie () (*in module dipper.sources.ClinVar*), 36

**F**

Family (*class in dipper.models.Family*), 25  
Feature (*class in dipper.models.GenomicFeature*), 25  
fetch () (*dipper.sources.AnimalQTLdb.AnimalQTLdb method*), 33  
fetch () (*dipper.sources.Bgee.Bgee method*), 34  
fetch () (*dipper.sources.BioGrid.BioGrid method*), 34  
fetch () (*dipper.sources.Coriell.Coriell method*), 38  
fetch () (*dipper.sources.CTD.CTD method*), 35  
fetch () (*dipper.sources.EBIGene2Phen.EBIGene2Phen method*), 39  
fetch () (*dipper.sources.Ensembl.Ensembl method*), 40  
fetch () (*dipper.sources.EOM.EOM method*), 39  
fetch () (*dipper.sources.FlyBase.FlyBase method*), 41  
fetch () (*dipper.sources.GeneOntology.GeneOntology method*), 42  
fetch () (*dipper.sources.GWASCatalog.GWASCatalog method*), 41  
fetch () (*dipper.sources.HPOAnnotations.HPOAnnotations method*), 43  
fetch () (*dipper.sources.IMPC.IMPC method*), 44  
fetch () (*dipper.sources.MGI.MGI method*), 45  
fetch () (*dipper.sources.MGISlim.MGISlim method*), 45  
fetch () (*dipper.sources.MMRRC.MMRRC method*), 46  
fetch () (*dipper.sources.Monarch.Monarch method*), 47  
fetch () (*dipper.sources.Monochrom.Monochrom method*), 48  
fetch () (*dipper.sources.MPD.MPD method*), 47  
fetch () (*dipper.sources.MyChem.MyChem method*), 49  
fetch () (*dipper.sources.MyDrug.MyDrug method*), 49  
fetch () (*dipper.sources.Orphanet.Orphanet method*), 50  
fetch () (*dipper.sources.Panther.Panther method*), 51  
fetch () (*dipper.sources.PostgreSQLSource.PostgreSQLSource method*), 51  
fetch () (*dipper.sources.Reactome.Reactome method*), 52  
fetch () (*dipper.sources.RGD.RGD method*), 52  
fetch () (*dipper.sources.SGD.SGD method*), 53  
fetch () (*dipper.sources.Source.Source method*), 54  
fetch () (*dipper.sources.StringDB.StringDB method*), 56  
fetch () (*dipper.sources.UCSCBands.UCSCBands method*), 57  
fetch () (*dipper.sources.UDP.UDP method*), 58  
fetch () (*dipper.sources.WormBase.WormBase method*), 59  
fetch () (*dipper.sources.Xenbase.Xenbase method*), 60  
fetch () (*dipper.sources.ZFIN.ZFIN method*), 61  
fetch () (*dipper.sources.ZFINSlim.ZFINSlim method*), 61  
fetch\_from\_mychem () (*dipper.sources.MyChem.MyChem method*), 49  
fetch\_from\_pgdb () (*dipper.sources.PostgreSQLSource.PostgreSQLSource method*), 51  
fetch\_from\_url () (*dipper.sources.Source.Source method*), 54  
fetch\_protein\_gene\_map () (*dipper.sources.Ensembl.Ensembl method*), 40  
fetch\_query\_from\_pgdb () (*dipper.sources.PostgreSQLSource.PostgreSQLSource method*), 51

fetch\_transgene\_genes\_from\_db() (dipper.sources.MGI.MGI method), 45  
fetch\_uniprot\_gene\_map() (dipper.sources.Ensembl.Ensembl method), 40  
fhandle (dipper.graph.RDFGraph.RDFGraph attribute), 14  
fhandle (dipper.graph.StreamedGraph.StreamedGraph attribute), 14  
fhandle (dipper.sources.ZFIN.ZFIN attribute), 61  
file\_len() (dipper.sources.Source.Source static method), 54  
files (dipper.sources.AnimalQTLdb.AnimalQTLdb attribute), 33  
files (dipper.sources.Bgee.Bgee attribute), 34  
files (dipper.sources.BioGrid.BioGrid attribute), 34  
files (dipper.sources.Coriell.Coriell attribute), 38  
files (dipper.sources.CTD.CTD attribute), 35  
files (dipper.sources.EBIGene2Phen.EBIGene2Phen attribute), 39  
files (dipper.sources.Ensembl.Ensembl attribute), 40  
files (dipper.sources.EOM.EOM attribute), 39  
files (dipper.sources.FlyBase.FlyBase attribute), 41  
files (dipper.sources.GeneOntology.GeneOntology attribute), 42  
files (dipper.sources.GWASCatalog.GWASCatalog attribute), 41  
files (dipper.sources.HPOAnnotations.HPOAnnotations attribute), 43  
files (dipper.sources.IMPC.IMPC attribute), 44  
files (dipper.sources.MMRC.MMRC attribute), 46  
files (dipper.sources.Monarch.Monarch attribute), 47  
files (dipper.sources.Monochrom.Monochrom attribute), 48  
files (dipper.sources.MPD.MPD attribute), 47  
files (dipper.sources.MyDrug.MyDrug attribute), 50  
files (dipper.sources.Orphanet.Orphanet attribute), 50  
files (dipper.sources.Panther.Panther attribute), 51  
files (dipper.sources.PostgreSQLSource.PostgreSQLSource attribute), 52  
files (dipper.sources.Reactome.Reactome attribute), 52  
files (dipper.sources.RGD.RGD attribute), 52  
files (dipper.sources.SGD.SGD attribute), 53  
files (dipper.sources.Source.Source attribute), 54  
files (dipper.sources.UCSCBands.UCSCBands attribute), 57  
files (dipper.sources.UDP.UDP attribute), 58  
files (dipper.sources.WormBase.WormBase attribute), 59  
files (dipper.sources.Xenbase.Xenbase attribute), 60  
files (dipper.sources.ZFIN.ZFIN attribute), 61  
files (dipper.sources.ZFINSlim.ZFINSlim attribute), 62  
FlyBase (class in dipper.sources.FlyBase), 40  
FLYFTP (dipper.sources.FlyBase.FlyBase attribute), 41  
format\_actions() (dipper.sources.MyChem.MyChem static method), 49  
fromRoman() (in module dipper.utils.romanplus), 64

## G

G2PAssoc (class in dipper.models.assoc.G2PAssoc), 17  
gaf\_columns (dipper.sources.GeneOntology.GeneOntology attribute), 42  
Gene (class in dipper.models.ClinVarRecord), 19  
gene\_info\_columns (dipper.sources.AnimalQTLdb.AnimalQTLdb attribute), 33  
GENEINFO (dipper.sources.AnimalQTLdb.AnimalQTLdb attribute), 33  
GeneOntology (class in dipper.sources.GeneOntology), 42  
Genotype (class in dipper.models.ClinVarRecord), 19  
Genotype (class in dipper.models.Genotype), 26  
Genovar (class in dipper.models.ClinVarRecord), 19  
get() (in module dipper.curie\_map), 64  
get\_association\_id() (dipper.models.assoc.Association.Assoc method), 15  
get\_base() (dipper.utils.CurieUtil.CurieUtil method), 62  
get\_base() (in module dipper.curie\_map), 64  
get\_common\_files() (dipper.sources.HPOAnnotations.HPOAnnotations method), 43  
get\_config() (in module dipper.config), 64  
get\_curie() (dipper.utils.CurieUtil.CurieUtil method), 62  
get\_curie\_prefix() (dipper.utils.CurieUtil.CurieUtil method), 62  
get\_drug\_record() (dipper.sources.MyChem.MyChem static method), 49  
get\_file\_md5() (dipper.sources.Source.Source static method), 54  
get\_files() (dipper.sources.Source.Source method), 54  
get\_graph() (dipper.models.Dataset.Dataset method), 22  
get\_hgnc\_id\_from\_symbol() (dipper.utils.DipperUtil.DipperUtil static method), 62  
get\_homologene\_by\_gene\_num() (dipper.utils.DipperUtil.DipperUtil static method), 62  
get\_inchikeys() (dipper.sources.MyChem.MyChem static method), 49

```

get_license()      (dipper.models.Dataset.Dataset     getTestSuite()      (dipper.sources.Source.Source
    method), 22           method), 54
get_local_file_size()      (dip-          getTestSuite()      (dip-
    per.sources.Source.Source static method),   per.sources.UCSCBands.UCSCBands method),
    54                      57
get_ncbi_taxon_num_by_label()      (dip-          gff_columns (dipper.sources.AnimalQTLdb.AnimalQTLdb
    per.utils.DipperUtil.DipperUtil static method),   attribute), 33
    62
get_orthology_evidence_code()      (dip-          GHRAW (dipper.sources.EOM.EOM attribute), 39
    per.sources.ZFIN.ZFIN static method), 61
get_orthology_sources_from_zebrafishmine()      (dip-          GITDIP (dipper.sources.AnimalQTLdb.AnimalQTLdb
    per.sources.ZFIN.ZFIN method), 61           attribute), 33
get_properties_from_graph()      (dip-          globaltcid (dipper.graph.RDFGraph.RDFGraph
    per.utils.GraphUtils.GraphUtils static method),   attribute), 14
    63
get_remote_content_len()      (dip-          globalttt (dipper.graph.RDFGraph.RDFGraph
    per.sources.Source method), 54           attribute), 14
get_uniprot_entrez_id_map()      (dip-          globalttt (dipper.graph.StreamedGraph.StreamedGraph
    per.sources.GeneOntology.GeneOntology           attribute), 14
    method), 42
get_uri()      (dipper.utils.CurieUtil.CurieUtil method), 41
    62
getChrPartTypeByNotation()      (in module dip-          Graph (class in dipper.graph.Graph), 13
    per.sources.Monochrom), 49
getTestSuite()      (dip-          GraphUtils (class in dipper.utils.GraphUtils), 63
    per.sources.AnimalQTLdb.AnimalQTLdb           GWASCatalog (class in dipper.sources.GWASCatalog),
    method), 33           41
getTestSuite()      (dipper.sources.BioGrid.BioGrid          GWASFILE (dipper.sources.GWASCatalog.GWASCatalog
    method), 34           attribute), 41
getTestSuite()      (dipper.sources.Coriell.Coriell          GWASFTP (dipper.sources.GWASCatalog.GWASCatalog
    method), 38           attribute), 41
getTestSuite()      (dipper.sources.CTD.CTD method), 35
getTestSuite()      (dipper.sources.Ensembl.Ensembl          H
    method), 40
getTestSuite()      (dipper.sources.EOM.EOM method), 39
getTestSuite()      (dipper.sources.GeneOntology.GeneOntology          hash_id() (dipper.models.Dataset.Dataset static
    method), 42           method), 22
getTestSuite()      (dip-          hash_id() (dipper.sources.Source.Source static
    per.sources.HPOAnnotations.HPOAnnotations           method), 54
    method), 43
getTestSuite()      (dipper.sources.IMPC.IMPC          HGGP (dipper.sources.UCSCBands.UCSCBands attribute),
    method), 44           57
getTestSuite()      (dipper.sources.MMRRC.MMRRC          HPOAnnotations (class in dipper.sources.HPOAnnotations), 42
    method), 46
getTestSuite()      (dip-          I
    per.sources.Monochrom.Monochrom method), 48
getTestSuite()      (dipper.sources.MPD.MPD          IMPC (class in dipper.sources.IMPC), 43
    method), 47
getTestSuite()      (dipper.sources.Panther.Panther          InteractionAssoc (class in dipper.models.assoc.InteractionAssoc), 18
    method), 51

```

## H

```

hash_id()      (dipper.models.Dataset.Dataset static
    method), 22
hash_id()      (dipper.sources.Source.Source static
    method), 54
HGGP          (dipper.sources.UCSCBands.UCSCBands attribute), 57
HPOAnnotations (class in dipper.sources.HPOAnnotations), 42

```

## I

```

IMPC          (class in dipper.sources.IMPC), 43
InteractionAssoc (class in dipper.models.assoc.InteractionAssoc), 18
is_id_in_mondo()      (dip-
    per.utils.DipperUtil.DipperUtil static method),
    62
is_literal()      (in module dipper.sources.ClinVar), 36

```

## K

```

key (dipper.models.BiolinkVocabulary.BioLinkVocabulary attribute), 18

```

## L

```

load_local_translationtable()      (dip-
    per.sources.Source method), 55

```

### M

make\_apo\_map() (*dipper.sources.SGD.SGD static method*), 53  
 make\_association() (*dipper.sources.RGD.RGD method*), 52  
 make\_association() (*dipper.sources.SGD.SGD method*), 53  
 make\_association\_id() (*dipper.models.assoc.Association.Assoc static method*), 15  
 make\_biolink\_category\_triple() (*in module dipper.sources.ClinVar*), 36  
 make\_c2p\_assoc\_id() (*dipper.models.assoc.Chem2DiseaseAssoc.Chem2DiseaseAssoc static method*), 16  
 make\_d2p\_id() (*dipper.models.assoc.D2PAssoc.D2PAssoc method*), 17  
 make\_experimental\_model\_with\_genotype() (*dipper.models.Genotype.Genotype method*), 29  
 make\_g2p\_id() (*dipper.models.assoc.G2PAssoc.G2PAssoc method*), 17  
 make\_id() (*dipper.models.Dataset.Dataset static method*), 22  
 make\_id() (*dipper.sources.Source.Source static method*), 55  
 make\_parent\_bands() (*dipper.sources.Monochrom.Monochrom method*), 48  
 make\_reagent\_targeted\_gene\_id() (*dipper.sources.WormBase.WormBase static method*), 59  
 make\_spo() (*in module dipper.sources.ClinVar*), 36  
 make\_targeted\_gene\_id() (*dipper.sources.ZFIN.ZFIN static method*), 61  
 make\_triples() (*dipper.sources.MyChem.MyChem method*), 49  
 make\_variant\_locus\_label() (*dipper.models.Genotype.Genotype static method*), 29  
 make\_vslc\_label() (*dipper.models.Genotype.Genotype method*), 29  
 makeChromID() (*in module dipper.models.GenomicFeature*), 26  
 makeChromLabel() (*in module dipper.models.GenomicFeature*), 26  
 makeGenomeID() (*dipper.models.Genotype.Genotype static method*), 29  
 makeLeader() (*dipper.models.Model.Model method*), 31  
 map\_files (*dipper.sources.EBIGene2Phen.EBIGene2Phen attribute*), 39  
 map\_files (*dipper.sources.UDP.UDP attribute*), 58  
 map\_type\_of\_region() (*dipper.sources.Monochrom.Monochrom method*), 48  
 mgd\_agent\_id (*dipper.sources.MPD.MPD attribute*), 47  
 mgd\_agent\_label (*dipper.sources.MPD.MPD attribute*), 47  
 mgd\_agent\_type (*dipper.sources.MPD.MPD attribute*), 47  
 MGI (*class in dipper.sources.MGI*), 45  
 MGISlim (*class in dipper.sources.MGISlim*), 45  
 MMRRC (*class in dipper.sources.MMRRC*), 46  
 Model (*class in dipper.models.Model*), 29  
 Monarch (*class in dipper.sources.Monarch*), 47  
 Monochrom (*class in dipper.sources.Monochrom*), 47  
 MPD (*class in dipper.sources.MPD*), 46  
 MPDDL (*dipper.sources.MPD.MPD attribute*), 47  
 MY\_DRUG\_API (*dipper.sources.MyDrug.MyDrug attribute*), 49  
 MyChem (*class in dipper.sources.MyChem*), 49  
 MyDrug (*class in dipper.sources.MyDrug*), 49

### N

namespaces (*dipper.sources.Source.Source attribute*), 55

### O

open\_and\_parse\_yaml() (*dipper.sources.Source.Source static method*), 55  
 Orphanet (*class in dipper.sources.Orphanet*), 50  
 OrthologyAssoc (*class in dipper.models.assoc.OrthologyAssoc*), 18

### P

Panther (*class in dipper.sources.Panther*), 50  
 panther\_format (*dipper.sources.Panther.Panther attribute*), 51  
 parse() (*dipper.sources.AnimalQTLdb.AnimalQTLdb method*), 33  
 parse() (*dipper.sources.Bgee.Bgee method*), 34  
 parse() (*dipper.sources.BioGrid.BioGrid method*), 35  
 parse() (*dipper.sources.Coriell.Coriell method*), 38  
 parse() (*dipper.sources.CTD.CTD method*), 35  
 parse() (*dipper.sources.EBIGene2Phen.EBIGene2Phen method*), 39  
 parse() (*dipper.sources.Ensembl.Ensembl method*), 40  
 parse() (*dipper.sources.EOM.EOM method*), 39  
 parse() (*dipper.sources.FlyBase.FlyBase method*), 41  
 parse() (*dipper.sources.GeneOntology.GeneOntology method*), 42

```

parse() (dipper.sources.GWASCatalog.GWASCatalog process_allele_phenotype() (dip-
    method), 41 per.sources.WormBase.WormBase method),
parse() (dipper.sources.HPOAnnotations.HPOAnnotations process_catalog() (dip-
    method), 43 per.sources.GWASCatalog.GWASCatalog method),
parse() (dipper.sources.IMPC.IMPC method), 44 process_common_disease_file() (dip-
parse() (dipper.sources.MGI.MGI method), 45 per.sources.HPOAnnotations.HPOAnnotations method),
process() (dipper.sources.MGISlim.MGISlim method), 45 process_disease_association() (dip-
    45 per.sources.WormBase.WormBase method),
parse() (dipper.sources.MMRRC.MMRRC method), 46 process_feature_loc() (dip-
    46 per.sources.WormBase.WormBase method),
parse() (dipper.sources.Monarch.Monarch method), 47 process_fish() (dipper.sources.ZFIN.ZFIN
    47 method), 59
parse() (dipper.sources.Monochrom.Monochrom process_fish_disease_models() (dip-
    method), 48 per.sources.ZFIN.ZFIN method),
parse() (dipper.sources.MPD.MPD method), 47 process_gaf() (dip-
parse() (dipper.sources.MyChem.MyChem method), 49 per.sources.GeneOntology.GeneOntology
    49 method),
parse() (dipper.sources.MyDrug.MyDrug method), 50 process_gene_desc() (dip-
parse() (dipper.sources.Orphanet.Orphanet method), 50 per.sources.WormBase.WormBase
    50 method),
process() (dipper.sources.Panther.Panther method), 51 process_gene_ids() (dip-
parse() (dipper.sources.PostgreSQLSource.PostgreSQLSource process_gene_interaction() (dip-
    method), 52 per.sources.WormBase.WormBase method),
parse() (dipper.sources.Reactome.Reactome method), 53 process_measure_set() (in module dip-
    53 per.sources.ClinVar), 36
parse() (dipper.sources.RGD.RGD method), 52
parse() (dipper.sources.SGD.SGD method), 53
parse() (dipper.sources.Source.Source method), 55
parse() (dipper.sources.StringDB.StringDB method), 57
parse() (dipper.sources.UCSCBands.UCSCBands process_mgi_note_allele_view() (dip-
    method), 58 per.sources.MGI.MGI method),
parse() (dipper.sources.UDP.UDP method), 58 process_mgi_relationship_transgene_genes() (dip-
parse() (dipper.sources.WormBase.WormBase per.sources.MGI.MGI method),
    59 method), 61
parse() (dipper.sources.Xenbase.Xenbase method), 60 process_omia_phenotypes() (dip-
parse() (dipper.sources.ZFIN.ZFIN method), 61 per.sources.Monarch.Monarch
    61 method),
parse() (dipper.sources.ZFINSlim.ZFINSlim method), 62 process_orthology_evidence() (dip-
    62 per.sources.ZFIN.ZFIN method),
parse() (in module dipper.sources.ClinVar), 36 process_pub_xrefs() (dip-
parse_checksum_file() (dip- per.sources.WormBase.WormBase
    per.sources.IMPC.IMPC method), 44 method), 60
parse_mapping_file() (dip- process_rnai_phenotypes() (dip-
    per.sources.Source.Source static method), 55 per.sources.WormBase.WormBase
    55 method),
Pathway (class in dipper.models.Pathway), 31 process_xml_table() (dip-
PNTHDL (dipper.sources.Panther.Panther attribute), 51 per.sources.Source method), 55
PostgreSQLSource (class in dip- Provenance (class in dipper.models.Provenance), 32
    per.sources.PostgreSQLSource), 51 Q
prefix_exists() (dipper.utils.CurieUtil.CurieUtil qt1_columns (dipper.sources.AnimalQTLdb.AnimalQTLdb
    method), 62
process_all_common_disease_files() (dip-
    per.sources.HPOAnnotations.HPOAnnotations
    method), 43

```

*attribute), 33*

*queries (dipper.sources.FlyBase.FlyBase attribute), 41*

**R**

*rdf2dot () (in module dipper.utils.rdf2dot), 64*

*RDFGraph (class in dipper.graph.RDFGraph), 13*

*Reactome (class in dipper.sources.Reactome), 52*

*REACTOME\_BASE (dipper.sources.Reactome.Reactome attribute), 52*

*record\_to\_triples () (in module dipper.sources.ClinVar), 36*

*Reference (class in dipper.models.Reference), 32*

*remove\_backslash\_r () (dipper.sources.Source.Source static method), 55*

*remove\_control\_characters () (dipper.utils.DipperUtil.DipperUtil static method), 63*

*remove\_ontology\_axioms () (dipper.utils.TestTools.TestTools static method), 63*

*resolve () (dipper.sources.Source.Source method), 56*

*resolve () (in module dipper.sources.ClinVar), 36*

*resources (dipper.sources.EOM.EOM attribute), 40*

*resources (dipper.sources.MGI.MGI attribute), 45*

*return\_target\_list () (dipper.sources.MyChem.MyChem static method), 49*

*RGD (class in dipper.sources.RGD), 52*

*RGD\_BASE (dipper.sources.RGD.RGD attribute), 52*

**S**

*scv\_link () (in module dipper.sources.ClinVar), 37*

*serialize () (dipper.graph.Graph.Graph method), 13*

*serialize () (dipper.graph.RDFGraph.RDFGraph method), 14*

*serialize () (dipper.graph.StreamedGraph.StreamedGraph method), 14*

*set\_association\_id () (dipper.models.assoc.Association.Assoc method), 15*

*set\_association\_id () (dipper.models.assoc.Chem2DiseaseAssoc.Chem2DiseaseAssoc method), 16*

*set\_association\_id () (dipper.models.assoc.D2PAssoc.D2PAssoc method), 17*

*set\_association\_id () (dipper.models.assoc.G2PAssoc.G2PAssoc method), 17*

*set\_citation () (dipper.models.Dataset.Dataset method), 22*

*set\_description () (dipper.models.assoc.Association.Assoc method), 16*

*set\_environment () (dipper.models.assoc.G2PAssoc.G2PAssoc method), 17*

*set\_ingest\_source () (dipper.models.Dataset.Dataset method), 22*

*set\_ingest\_source\_file\_version\_date () (dipper.models.Dataset.Dataset method), 22*

*set\_ingest\_source\_file\_version\_num () (dipper.models.Dataset.Dataset method), 23*

*set\_ingest\_source\_file\_version\_retrieved\_on () (dipper.models.Dataset.Dataset method), 23*

*set\_object () (dipper.models.assoc.Association.Assoc method), 16*

*set\_relationship () (dipper.models.assoc.Association.Assoc method), 16*

*set\_score () (dipper.models.assoc.Association.Assoc method), 16*

*set\_stage () (dipper.models.assoc.G2PAssoc.G2PAssoc method), 17*

*set\_subject () (dipper.models.assoc.Association.Assoc method), 16*

*setAuthorList () (dipper.models.Reference.Reference method), 32*

*setShortCitation () (dipper.models.Reference.Reference method), 32*

*settestmode () (dipper.sources.Source.Source method), 56*

*settestonly () (dipper.sources.Source.Source method), 56*

*setTitle () (dipper.models.Reference.Reference method), 32*

*setType () (dipper.models.Reference.Reference method), 32*

*setYear () (dipper.models.Reference.Reference method), 32*

*SGD (class in dipper.sources.SGD), 53*

*SGD\_BASE (dipper.sources.SGD.SGD attribute), 53*

*skolemizeBlankNode () (dipper.graph.Graph.Graph method), 13*

*skolemizeBlankNode () (dipper.graph.RDFGraph.RDFGraph method), 14*

*skolemizeBlankNode () (dipper.graph.StreamedGraph.StreamedGraph method), 14*

*small\_files (dipper.sources.HPOAnnotations.HPOAnnotations attribute), 43*

*Source (class in dipper.sources.Source), 53*

species (*dipper.sources.WormBase.WormBase attribute*), 60  
 StreamedGraph (*class in dipper.graph.StreamedGraph*), 14  
 StringDB (*class in dipper.sources.StringDB*), 56

**T**

tables (*dipper.sources.EOM.EOM attribute*), 40  
 tables (*dipper.sources.MGI.MGI attribute*), 45  
 terms (*dipper.models.BiolinkVocabulary.BioLinkVocabulary attribute*), 18  
 test\_graph\_equality () (*dipper.utils.TestTools.TestTools static method*), 63  
 test\_ids (*dipper.sources.AnimalQTLdb.AnimalQTLdb attribute*), 33  
 test\_ids (*dipper.sources.MMRRC.MMRRC attribute*), 46  
 test\_ids (*dipper.sources.MPD.MPD attribute*), 47  
 test\_ids (*dipper.sources.ZFIN.ZFIN attribute*), 61  
 test\_lines (*dipper.sources.Coriell.Coriell attribute*), 38  
 TestUtils (*class in dipper.utils.TestTools*), 63  
 toRoman () (*in module dipper.utils.romanplus*), 64  
 trait\_mapping\_columns (*dipper.sources.AnimalQTLdb.AnimalQTLdb attribute*), 33

**U**

UCSCBands (*class in dipper.sources.UCSCBands*), 57  
 UDP (*class in dipper.sources.UDP*), 58  
 UDP\_SERVER (*dipper.sources.UDP.UDP attribute*), 58  
 unknown\_taxa (*dipper.sources.MGI.MGI attribute*), 45  
 update\_wsnum\_in\_files () (*dipper.sources.WormBase.WormBase method*), 60

**V**

Variant (*class in dipper.models.ClinVarRecord*), 19

**W**

wbdev (*dipper.sources.WormBase.WormBase attribute*), 60  
 wbprod (*dipper.sources.WormBase.WormBase attribute*), 60  
 wbrel (*dipper.sources.WormBase.WormBase attribute*), 60  
 whoami () (*dipper.sources.Source.Source method*), 56  
 wont\_prefix (*dipper.sources.GeneOntology.GeneOntology attribute*), 42  
 WormBase (*class in dipper.sources.WormBase*), 59  
 write () (*dipper.sources.Source.Source method*), 56

at- write () (*dipper.utils.GraphUtils.GraphUtils static method*), 63  
 dip- write\_review\_status\_scores () (*in module dipper.sources.ClinVar*), 37  
 write\_spo () (*in module dipper.sources.ClinVar*), 37

**X**

Xenbase (*class in dipper.sources.Xenbase*), 60

**Y**

yaml\_file (*dipper.models.BiolinkVocabulary.BioLinkVocabulary attribute*), 18

**Z**

ZFIN (*class in dipper.sources.ZFIN*), 60  
 ZFINSlim (*class in dipper.sources.ZFINSlim*), 61