
pypath Documentation

Release 0.7.97

Dénes Türei

Oct 22, 2018

CONTENTS:

1	Webservice	3
1.1	Mouse and rat	3
1.2	Examples	3
2	Can I use OmniPath in R?	5
3	Installation	7
3.1	Linux	7
3.2	igraph C library, cairo and pycairo	7
3.3	Directly from git	7
3.4	With pip	7
3.5	Build source distribution	7
3.6	Mac OS X	8
3.7	Microsoft Windows	8
4	Release History	11
4.1	0.1.0:	11
4.2	0.2.0:	11
4.3	0.3.0:	11
4.4	0.4.0:	11
4.5	0.5.0:	11
4.6	0.7.74:	12
4.7	Upcoming:	12
5	Features	13
5.1	ID conversion	13
5.2	Pathways	13
5.3	Structural features	13
5.4	Sequences	13
5.5	Tissue expression	14
5.6	Functional annotations	14
5.7	Drug compounds	14
5.8	Technical	14

note `pypath` supports both Python 2.7 and Python 3.6+. In the beginning, `pypath` has been developed only for Python 2.7. Then the code have been adjusted to Py3 however we can not guarantee no incompatibilities remained. If you find any method does not work please submit an issue on github. For few years I develop and test `pypath` in Python 3. Therefore this is the better supported Python variant.

contributions turei.denes@gmail.com

documentation <http://pypath.omnipathdb.org/>

issues <https://github.com/saezlab/pypath/issues>

pypath is a Python package built around `igraph` to work with molecular network representations e.g. protein, miRNA and drug compound interaction networks.

WEBSERVICE

New webservice from 14 June 2018: the queries slightly changed, have been largely extended. See the examples below.

One instance of the `pypath` webservice runs at the domain <http://omnipathdb.org/>, serving not only the OmniPath data but other datasets: TF-target interactions from TF Regulons, a large collection additional enzyme-substrate interactions, and literature curated miRNA-mRNA interactions combined from 4 databases. The webservice implements a very simple REST style API, you can make requests by HTTP protocol (browser, `wget`, `curl` or whatever).

The webservice currently recognizes 3 types of queries: `interactions`, `ptms` and `info`. The query types `resources`, `network` and `about` have not been implemented yet in the new webservice.

1.1 Mouse and rat

Except the miRNA interactions all interactions are available for human, mouse and rat. The rodent data has been translated from human using the NCBI Homologene database. Many human proteins have no known homolog in rodents hence rodent datasets are smaller than their human counterparts. Note, if you work with mouse omics data you might do better to translate your dataset to human (for example using the `pypath.homology` module) and use human interaction data.

1.2 Examples

A request without any parameter, gives some basic numbers about the actual loaded dataset:

<http://omnipathdb.org>

The `info` returns a HTML page with comprehensive information about the resources:

<http://omnipathdb.org/info>

The `interactions` query accepts some parameters and returns interactions in tabular format. This example returns all interactions of EGFR (P00533), with sources and references listed.

<http://omnipathdb.org/interactions/?partners=P00533&fields=sources,references>

By default only the OmniPath dataset used, to query the TF Regulons or add the extra enzyme-substrate interactions you need to set additional parameters. For example to query the transcriptional regulators of EGFR:

<http://omnipathdb.org/interactions/?targets=EGFR&types=TF>

The TF Regulons database assigns confidence levels to the interactions. You might want to select only the highest confidence, A category:

http://omnipathdb.org/interactions/?targets=EGFR&types=TF&tfregulons_levels=A

Show the transcriptional targets of Smad2 homology translated to rat including the confidence levels from TF Regulations:

```
http://omnipathdb.org/interactions/?genesymbols=1&fields=type,ncbi_tax_id,tfregulons_level&organisms=10116&sources=Smad2&types=TF
```

Query interactions from PhosphoNetworks which is part of the *kinaseextra* dataset:

```
http://omnipathdb.org/interactions/?genesymbols=1&fields=sources&databases=PhosphoNetworks&datasets=kinaseextra
```

Get the interactions from Signor, SPIKE and Signalink3:

```
http://omnipathdb.org/interactions/?genesymbols=1&fields=sources,references&databases=Signor,SPIKE,Signalink3
```

All interactions of MAP1LC3B:

```
http://omnipathdb.org/interactions/?genesymbols=1&partners=MAP1LC3B
```

By default `partners` queries the interaction where either the source or the target is among the partners. If you set the `source_target` parameter to AND both the source and the target must be in the queried set:

```
http://omnipathdb.org/interactions/?genesymbols=1&fields=sources,references&sources=ATG3,ATG7,ATG4B,SQSTM1&targets=MAP1LC3B,MAP1LC3A,MAP1LC3C,Q9H0R8,GABARAP,GABARAPL2&source_target=AND
```

As you see above you can use UniProt IDs and Gene Symbols in the queries and also mix them. Get the miRNA regulating NOTCH1:

```
http://omnipathdb.org/interactions/?genesymbols=1&fields=sources,references&datasets=mirnatarget&targets=NOTCH1
```

Note: with the exception of mandatory fields and genesymbols, the columns appear exactly in the order you provided in your query.

Another query type available is `ptms` which provides enzyme-substrate interactions. It is very similar to the interactions:

```
http://omnipathdb.org/ptms?genesymbols=1&fields=sources,references,isoforms&enzymes=FYN
```

Is there any ubiquitination reaction?

```
http://omnipathdb.org/ptms?genesymbols=1&fields=sources,references&types=ubiquitination
```

And acetylation in mouse?

```
http://omnipathdb.org/ptms?genesymbols=1&fields=sources,references&types=acetylation&organisms=10090
```

Rat interactions, both directly from rat and homology translated from human, from the PhosphoSite database:

```
http://omnipathdb.org/ptms?genesymbols=1&fields=sources,references&organisms=10116&databases=PhosphoSite,PhosphoSite_noref
```


CAN I USE OMNIPATH IN R?

You can download the data from the webservice and load into R. Look here for an example:

https://github.com/saezlab/pypath/tree/master/r_import

INSTALLATION

3.1 Linux

In almost any up-to-date Linux distribution the dependencies of **pypath** are built-in, or provided by the distributors. You only need to install a couple of things in your package manager (cairo, py(2)cairo, igraph, python(2)-igraph, graphviz, pygraphviz), and after install **pypath** by *pip* (see below). If any module still missing, you can install them the usual way by *pip* or your package manager.

3.2 igraph C library, cairo and pycairo

python(2)-igraph is a Python interface to use the igraph C library. The C library must be installed. The same goes for *cairo*, *py(2)cairo* and *graphviz*.

3.3 Directly from git

```
pip install git+https://github.com/saezlab/pypath.git
```

3.4 With pip

Download the package from /dist, and install with pip:

```
pip install pypath-x.y.z.tar.gz
```

3.5 Build source distribution

Clone the git repo, and run setup.py:

```
python setup.py sdist
```

3.6 Mac OS X

On OS X installation is not straightforward primarily because cairo needs to be compiled from source. We provide 2 scripts here: the **mac-install-brew.sh** installs everything with HomeBrew, and **mac-install-conda.sh** installs from Anaconda distribution. With these scripts installation of igraph, cairo and graphviz goes smoothly most of the time, and options are available for omitting the 2 latter. To know more see the description in the script header. There is a third script **mac-install-source.sh** which compiles everything from source and presumes only Python 2.7 and Xcode installed. We do not recommend this as it is time consuming and troubleshooting requires expertise.

3.6.1 Troubleshooting

- no module named ... when you try to load a module in Python. Did the installation of the module run without error? Try to run again the specific part from the mac install shell script to see if any error comes up. Is the path where the module has been installed in your `$PYTHONPATH`? Try `echo $PYTHONPATH` to see the current paths. Add your local install directories if those are not there, e.g. `export PYTHONPATH="/Users/me/local/python2.7/site-packages:$PYTHONPATH"`. If it works afterwards, don't forget to append these export path statements to your `~/.bash_profile`, so these will be set every time you launch a new shell.
- `pkgconfig` not found. Check if the `$PKG_CONFIG_PATH` variable is set correctly, and pointing on a directory where `pkgconfig` really can be found.
- Error while trying to install `py(2)cairo` by `pip`. `py(2)cairo` could not be installed by `pip`, but only by `waf`. Please set the `$PKG_CONFIG_PATH` before. See **mac-install-source.sh** on how to install with `waf`.
- Error at `pygraphviz` build: `graphviz/cgraph.h` file not found. This is because the directory of `graphviz` detected wrong by `pkgconfig`. See **mac-install-source.sh** how to set include dirs and library dirs by `--global-option` parameters.
- Can not install `bioservices`, because installation of `jurko-suds` fails. Ok, this fails because `pip` is not able to install the recent version of `setuptools`, because a very old version present in the system path. The development version of `jurko-suds` does not require `setuptools`, so you can install it directly from git as it is done in **mac-install-source.sh**.
- In **Anaconda**, `pypath` can be imported, but the modules and classes are missing. Apparently Anaconda has some built-in stuff called `pypath`. This has nothing to do with this module. Please be aware that Anaconda installs a completely separated Python distribution, and does not detect modules in the main Python installation. You need to install all modules within Anaconda's directory. **mac-install-conda.sh** does exactly this. If you still experience issues, please contact us.

3.7 Microsoft Windows

Not many people have used `pypath` on Microsoft computers so far. Please share your experiences and contact us if you encounter any issue. We appreciate your feedback, and it would be nice to have better support for other computer systems.

3.7.1 With Anaconda

The same workflow like you see in `mac-install-conda.sh` should work for Anaconda on Windows. The only problem you certainly will encounter is that not all the channels have packages for all platforms. If certain channel provides no package for Windows, or for your Python version, you just need to find an other one. For this, do a search:

```
anaconda search -t conda <package name>
```

For example, if you search for *pycairo*, you will find out that *vgauther* provides it for osx-64, but only for Python 3.4, while *richlewis* provides also for Python 3.5. And for win-64 platform, there is the channel of *KristanAmstrong*. Go along all the commands in `mac-install-conda.sh`, and modify the channel if necessary, until all packages install successfully.

3.7.2 With other Python distributions

Here the basic principles are the same as everywhere: first try to install all external dependencies, after *pip* install should work. On Windows certain packages can not be installed by compiled from source by *pip*, instead the easiest to install them precompiled. These are in our case *fisher*, *lxml*, *numpy (mkl version)*, *pycairo*, *igraph*, *pygraphviz*, *scipy* and *statsmodels*. The precompiled packages are available here: <http://www.lfd.uci.edu/~gohlke/pythonlibs/>. We tested the setup with Python 3.4.3 and Python 2.7.11. The former should just work fine, while with the latter we have issues to be resolved.

3.7.3 Known issues

- “No module fabric available.” – or *pysftp* missing: this is not important, only certain data download methods rely on these modules, but likely you won’t call those at all.
- Progress indicator floods terminal: sorry about that, will be fixed soon.
- Encoding related exceptions in Python2: these might occur at some points in the module, please send the traceback if you encounter one, and we will fix as soon as possible.

Special thanks to Jorge Ferreira for testing pypath on Windows!

RELEASE HISTORY

Main improvements in the past releases:

4.1 0.1.0:

- First release of pypath, for initial testing.

4.2 0.2.0:

- Lots of small improvements in almost every module
- Networks can be read from local files, remote files, lists or provided by any function
- Almost all redistributed data have been removed, every source downloaded from the original provider.

4.3 0.3.0:

- First version with partial Python 3 support.

4.4 0.4.0:

- **pyreact** module with **BioPaxReader** and **PyReact** classes added
- Process description databases, BioPax and PathwayCommons SIF conversion rules are supported
- Format definitions for 6 process description databases included.

4.5 0.5.0:

- Many classes have been added to the **plot** module
- All figures and tables in the manuscript can be generated automatically
- This is supported by a new module, **analysis**, which implements a generic workflow in its **Workflow** class.

4.6 0.7.74:

- **homology** module: finds the homologs of proteins using the NCBI

Homologene database and the homologs of PTM sites using UniProt sequences and PhosphoSitePlus homology table
* **ptm** module: fully integrated way of processing enzyme-substrate interactions from many databases and their translation by homology to other species
* **export** module: creates `pandas.DataFrame` or exports the network into tabular file
* New webservice
* TF Regulons database included and provides much more comprehensive transcriptional regulation resources, including literature curated, in silico predicted, ChIP-Seq and expression pattern based approaches
* Many network resources added, including miRNA-mRNA and TF-miRNA interactions

4.7 Upcoming:

- New, more flexible network reader class
- Full support for multi-species molecular interaction networks

(e.g. pathogene-host)
* Better support for not protein only molecular interaction networks (metabolites, drug compounds, RNA)
* ChEMBL webservice interface, interface for PubChem and eventually for DrugBank
* Silent mode: a way to suppress messages and progress bars

FEATURES

The primary aim of **pypath** is to build up networks from multiple sources on one **igraph** object. **pypath** handles ambiguous ID conversion, reads custom edge and node attributes from text files and **MySQL**.

Submodules perform various features, e.g. graph visualization, working with **rug** compound data, searching drug targets and compounds in **ChEMBL**.

5.1 ID conversion

The ID conversion module `mapping` can be used independently. It has the feature to translate secondary UniProt IDs to primaries, and Trembl IDs to SwissProt, using primary Gene Symbols to find the connections. This module automatically loads and stores the necessary conversion tables. Many tables are predefined, such as all the IDs in **UniProt mapping service**, while users are able to load any table from **file** or **MySQL**, using the classes provided in the module `input_formats`.

5.2 Pathways

pypath includes data and predefined format descriptions for more than 25 high quality, literature curated databases. The input formats are defined in the `data_formats` module. For some resources data downloaded on the fly, where it is not possible, data is redistributed with the module. Descriptions and comprehensive information about the resources is available in the `descriptions` module.

5.3 Structural features

One of the modules called `intera` provides many classes for representing structures and mechanisms behind protein interactions. These are `Residue` (optionally mutated), `Motif`, `Ptm`, `Domain`, `DomainMotif`, `DomainDomain` and `Interface`. All these classes have `__eq__()` methods to test equality between instances, and also `__contains__()` methods to look up easily if a residue is within a short motif or protein domain, or is the target residue of a PTM.

5.4 Sequences

The module `seq` contains a simple class for quick lookup any residue or segment in **UniProt** protein sequences while being aware of isoforms.

5.5 Tissue expression

For 3 protein expression databases there are functions and modules for downloading and combining the expression data with the network. These are the Human Protein Atlas, the ProteomicsDB and GIANT. The `giant` and `proteomicsdb` modules can be used also as stand alone Python clients for these resources.

5.6 Functional annotations

GSEA and **Gene Ontology** are two approaches for annotating genes and gene products, and enrichment analysis technics aims to use these annotations to highlight the biological functions a given set of genes is related to. Here the `enrich` module gives abstract classes to calculate enrichment statistics, while the `go` and the `gsea` modules give access to GO and GSEA data, and make it easy to count enrichment statistics for sets of genes.

5.7 Drug compounds

UniChem submodule provides an interface to effectively query the UniChem service, use connectivity search with custom settings, and translate SMILES to ChEMBL IDs with ChEMBL web service.

ChEMBL submodule queries directly your own ChEMBL MySQL instance, has the features to search targets and compounds from custom assay types and relationship types, to get activity values, binding domains, and action types. You need to download the ChEMBL MySQL dump, and load into your own server.

5.8 Technical

MySQL submodule helps to manage MySQL connections and track queries. It is able to run queries parallelly to optimize CPU and memory usage on the server, handling queues, and serve the result by server side or client side storage. The `chembl` and potentially the `mapping` modules rely on this `mysql` module.

The most important function in module `dataio` is a very flexible **download manager** built around `curl`. The function `dataio.curl()` accepts numerous arguments, tries to deal in a smart way with local **cache**, authentication, redirects, uncompression, character encodings, FTP and HTTP transactions, and many other stuff. Cache can grow to several GBs, and takes place in `./cache` by default. Please be aware of this, and use for example symlinks in case of using multiple working directories.

A simple **webservice** comes with this module: the `server` module based on `twisted.web.server` opens a custom port and serves plain text tables over HTTP with REST style querying.