



NodeSoftware Documentation

Document Information

Editors: T. Marquart
Authors: T. Marquart, S. Regandell, G. Rixon
Contributors: VAMDC WP7 working group
Type of document: software documentation
Status: release
Distribution: public
Work package: WP7
Version: 11.5r1
Document code:
Directory and file name:

Abstract: This document describes the functionality and the setup of the VAMDC NodeSoftware which implements the standards for a VAMDC service.

Version History

Version	Date	Modified By	Description of Change
0.1	16/12/2010	T. Marquart	first draft
0.2	27/01/2011	T. Marquart	version for implementation workshop
11.5	27/05/2011	T. Marquart	release together with standards
11.5r1	15/06/2011	T. Marquart	update to match software release 11.5r1

Disclaimer

The information in this document is subject to change without notice. Company or product names mentioned in this document may be trademarks or registered trademarks of their respective companies.

All rights reserved

The document is proprietary of the VAMDC consortium members. No copying or distributing, in any form or by any means, is allowed without the prior written agreement of the owner of the property rights.

This document reflects only the authors' view. The European Community is not liable for any use that may be made of the information contained herein.

Acknowledgements

VAMDC is funded under the "Combination of Collaborative Projects and Coordination and Support Actions" Funding Scheme of The Seventh Framework Program. Call topic: INFRA-2008-1.2.2 Scientific Data Infrastructure. Grant Agreement number: 239108.

CONTENTS

1	Introduction	2
1.1	About VAMDC	2
1.2	VAMDC nodes	2
1.3	A versatile implementation of VAMDC standards	2
2	Changelog	4
2.1	June 15, 2011	4
2.2	May 26, 2011	4
2.3	March 10, 2011	5
2.4	February 2011	6
3	The main concepts behind the implementation	7
3.1	The database	7
3.2	The data model(s)	7
3.3	The VAMDC dictionary	7
3.4	The registry	9
3.5	TAP services	9
3.6	The query language	9
3.7	The XSAMS schema	9
3.8	The generic XSAMS generator	10
3.9	The portal	10
4	Software prerequisites and installation	11
4.1	Quick start	11
4.2	Python plus some modules	11
4.3	Django	12
4.4	Database engine	12
4.5	Webserver	12
4.6	Git version control	12
4.7	The node software itself	12
4.8	Test your installation	12
5	Upgrading	14
5.1	NodeSoftware	14
5.2	Django	14
5.3	Everything else	14
6	Step by step guide to a new VAMDC node	15
6.1	The main directory of your node	16
6.2	Inside your node directory	16
6.3	The data model and the database	16
6.4	Using the XML generator	19
6.5	The query routine	20
6.6	The dictionaries	23

6.7	Testing the node	24
7	How to get your data into the database	25
7.1	Loading ascii data into the database	25
7.2	Preparing the input files	26
7.3	The mapping file	26
8	Deployment of your node	31
8.1	Gunicorn plus proxy	31
8.2	Deployment in Apache	32
8.3	Third party hosting	33
8.4	Logging	33
9	Obtain and use a Virtual Machine with the NodeSoftware	34
9.1	About VirtualBox	34
9.2	The virtual harddisk	34
9.3	Setting up the VM	34
9.4	Once inside the VM	34
10	Additional topics	37
10.1	Setting the related name of a field	37
10.2	Inserting custom XML into the generator	37
10.3	Collaborating with git and GitHub	37
10.4	Adding more views or apps to your node	40
10.5	The Django admin interface	41
10.6	Handling advanced queries	41
10.7	Using a custom model method for filling a Returnable	41
11	Known limitations	42
12	Bugs and Contact	43
12.1	Report a bug	43
12.2	Contact information	43
13	The Code	44
13.1	Source code documentation	44
13.2	The VAMDC-TAP service library	44
13.3	The import tool	48
	Python Module Index	51
	Index	52

This document covers the **release 11.5r1** of the NodeSoftware.

PDF-versions of this document are available at:

- http://www.vamdc.org/documents/NodeSoftwareDoc_v11.5r1.pdf (release version)
- <http://vamdc.tmy.se/doc/nodesoftware.pdf> (latest development version)

INTRODUCTION

1.1 About VAMDC

The Virtual Atomic and Molecular Data Center is a EU FP7 research infrastructure project and you can read all about it on <http://vamdc.eu/>

1.2 VAMDC nodes

A “node” within VAMDC is a data service that offers its data using the standards and protocols defined by the VAMDC. They are RESTful HTTP services, the specification of which can be found in the documentation for the VAMDC standards: <http://vamdc.org/documents/standards/>

The scope of this document is to serve as documentation for the reference implementation of such a service. The goal of this implementation is to serve as publishing tools for new data services, i.e. it is meant to be easily deployed at multiple nodes.

1.3 A versatile implementation of VAMDC standards

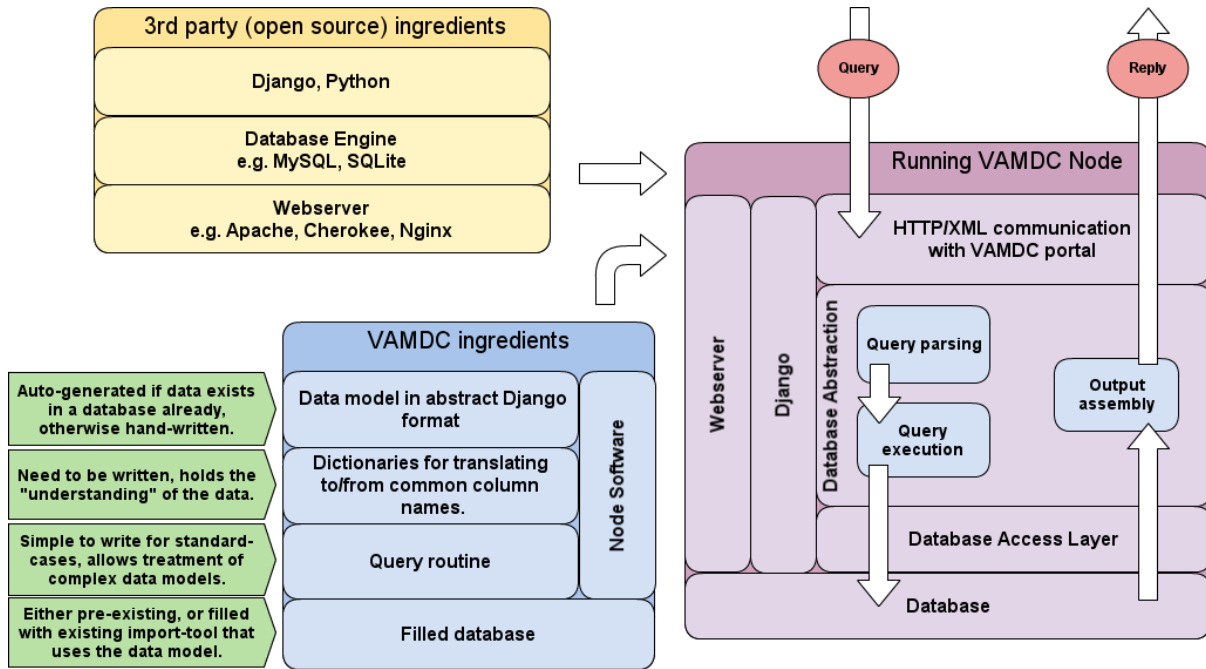
Principle design decisions that were made to arrive at this software package include

- *Open source.* No software licences need to be bought and the used software can be adapted if needed.
- The data must exist in a *relational database*. If this is not the case yet, a tool for creating it is provided.
- *Flexibility in the data structure.* The service should be able to be plugged on top of existing databases and therefore needs to cope with almost arbitrary internal data formats.
- *Re-usable code.* The implementation of the VAMDC standards and protocols themselves should not depend on the requirements of a specific node.

Since the last two points contradict each other in practice, there needs to be an intermediate layer of abstraction that hides the node-specific details like the database layout from the parts of code that are shared between nodes.

Our implementation of the VAMDC node software is therefore based on a framework called *Django* (which in turn is based on the programming language *Python*) that provides both the database abstraction layer and high level tools for implementing web services.

The ingredients for a VAMDC node based on this software package and its operation look schematically like this:



CHANGELOG

This chapter will be difficult to understand if you have not read the whole document before since terms are used that are introduced later. It is meant for returning readers, especially the maintainer of VAMDC nodes.

2.1 June 15, 2011

Version. This documentation has been updated to match the release of the NodeSoftware 11.5r1 which implements the VAMDC Standards release 11.5. NodeSoftware 11.5r1 supersedes and obsoletes version 11.5 (released May 26) and all nodes are encouraged to upgrade. This is mainly a bug-fix release and upgraded nodes will only have to do the two small changes mentioned below.

Example Queries. The way to define example queries in each node's `settings.py` has changed in order to allow several of them. They will be used for automated testing and are as of this version returned to the VAMDC registry. New example:

```
EXAMPLE_QUERIES = [\n    'SELECT ALL WHERE RadTransWavelength > 4000 AND RadTransWavelength < 4005',\n    'SELECT ALL WHERE AtomSymbol = "Fe"',\n]
```

CaselessDict. The import and use of *CaselessDict* in the nodes' `dictionaries.py` or `queryfunc.py` is not longer necessary and should be removed.

Limitations. A chapter on the limitations of the NodeSoftware has been added to the documentation: *Known limitations*

Dictionary. The NodeSoftware makes use of dictionary keywords that are not in the VAMDC Standards 11.5 but will be in the next Standards release (11.7). If you want to use the NodeSoftware's XML-generator for solids, particles or molecular quantum numbers, please see <http://dictionary.vamdc.org/dict/> for the new keywords.

Registration. The NodeSoftware now automatically reports its own version and the standards version it implements at *tap/capabilities*. You might want to make the VAMDC Registry re-read this information (click "Edit metadata" and "Update the registry entry").

Virtual Machine. The virtual machine has been updated to include Django 1.3 and NodeSoftware 11.5r1.

2.2 May 26, 2011

Version numbers. As of now, we introduce version numbers for both the standards (XSAMS, VAMDC-TAP, see separate documentation) and for their implementation in the NodeSoftware which is the concern of this document. Version numbers follow the format YY.MMrX where YY is for the year, MM the month, and X an increasing number for bugfix revisions that do not affect the usage of the NodeSoftware.

The most important changes from the perspective of a node-operator who wants to upgrade to this 11.5 release are:

Update to Django 1.3. The NodeSoftware now requires Django version 1.3 and node operators probably need to upgrade their installation of Django. See [Upgrading](#).

Email. Make sure you have set a correct email address in `settings.py`. It will be used to report critical errors to, including reports on what went wrong.

Logging. The capabilities to log debug and error-messages have been extended. See [Logging](#).

Example query. As soon as a node becomes operational, please add an example query to its `settings.py`. It will be used for automated testing. Example:

```
EXAMPLE_QUERY = 'SELECT ALL WHERE RadTransWavelength > 4000 AND RadTransWavelength < 4005'
```

Volume estimate. In order to allow the portal (and other queries to your node) to find out how big the resulting XML-output for a particular query will be, nodes should estimate this and relay it via the new HTTP-header `VAMDC-APPROX-SIZE`. The easiest way to do this is to run a test query, determine the outputs size (in MB) and divide it by the number of items (e.g. transitions, if these dominate your results). This number can then be used to estimate the size of any query, see the updated example at [The query routine](#).

Other Header changes. The header `VAMDC-COUNT-SPECIES` has been replaced by `VAMDC-COUNT-ATOMS` and `VAMDC-COUNT-MOLECULES`. See the standards documentation for the full definition.

Error handling in `urls.py`. The NodeSoftware has become more error-safe and tries to handle unexected input and “crashes” more gracefully. You need not care about this, except making sure that the following two lines are present at the end of the file `urls.py` in your node’s main directory:

```
handler500 = 'vamdctap.views.tapServerError'  
handler404 = 'vamdctap.views.tapNotFoundError'
```

Dictionary changes. Since the XSAMS-schema has changed, so have the dictionary keywords, especially in the Broadening-part of radiative transitions and the atomic quantum numbers. Also new keywords have been added for the bits that are newly implemented in the XML-generator.

Stricter format for accuracies. In compliance with XSAMS’ new way of defining a value’s accuracy, the keywords that are not explicitly given for *DataTypes* have become more. Any word *SomeKeyword* that is marked as a *DataType* in the dictionary allows for use of the following words as well: *SomeKeywordUnit*, *SomeKeywordRef*, *SomeKeywordComment*, *SomeKeywordMethod*, *SomeKeywordAccuracyCalibration*, *SomeKeywordAccuracyQuality*, *SomeKeywordAccuracySystematic*, *SomeKeywordAccuracySystematicConfidence*, *SomeKeywordAccuracySystematicRelative*, *SomeKeywordAccuracyStatistical*, *SomeKeywordAccuracyStatisticalConfidence*, *SomeKeywordAccuracyStatisticalRelative*, *SomeKeywordAccuracyStatLow*, *SomeKeywordAccuracyStatLowConfidence*, *SomeKeywordAccuracyStatLowRelative*, *SomeKeywordAccuracyStatHigh*, *SomeKeywordAccuracyStatHighConfidence*, *SomeKeywordAccuracyStatHighRelative*. See also the standards documentation.

Note: The last two points mean that you probably have to update your `dictionaries.py`.

2.3 March 10, 2011

The chapter [The main concepts behind the implementation](#) now has more detail on the XSAMS schema.

A large part of the XML/XSAMS generator has been rewritten, both to comply with the new version of the schema and in terms of its structure. In addition the keywords in the VAMDC dictionary have changed somewhat. This means that **you will probably need to update your query function and dictionaries when you update the NodeSoftware.**

[Step by step guide to a new VAMDC node](#) has been updated and extended accordingly.

A new version of the [Obtain and use a Virtual Machine with the NodeSoftware](#) has also been uploaded, containing the latest NodeSoftware and operating system.

2.4 February 2011

The deployment of nodes is now covered in more detail at *Deployment of your node*.

THE MAIN CONCEPTS BEHIND THE IMPLEMENTATION

The following is a glossary-like list that shortly touches upon various subjects that one should be aware of before setting up a new VAMDC node.

3.1 The database

As already mentioned in the *Introduction*, data needs to reside in a relational database in order to use the node software for a VAMDC node. This is what we mean by *database* in the following, in contrast to *data set* which means the data in any format or *data model*:

3.2 The data model(s)

The *data model* is a definition of the database layout in form of Python code where a *class* is defined for each table in the database and the members of the class are *fields* that correspond to the tables' columns. The data model also defines the connections between tables. For an existing database the data model can be automatically generated, otherwise it needs to be written for a new node (see *Step by step guide to a new VAMDC node* later) and will then be used to create the database.

Having this code representation of the database layout has many advantages, among these are:

- automatic (re-)creation of the database, independent of the engine
- no need to learn SQL
- easy queries
- additional features like easily traversing linked tables in both directions.

Note: Sometimes the singular *data model* refers to a single model (i.e. a table in the database) and sometimes the full set of models, describing the whole database layout.

3.3 The VAMDC dictionary

In order to facilitate automated communication, there is a need for a set of names that identify a certain type of data. Each name is unique and is uniquely associated with a description, a data type, a unit where applicable and a (non-mandatory) restriction.

For illustration, let's have a look at one entry of the dictionary:

Keyword	short descr	long description	data type	re-strict-ion	unit
Atom-Mass-Number	Atomic mass	Atomic mass in Daltons, which is the same as the unified mass units (1Da = 1u = 1.660 538 86 (28) e-27)	(Float Double)		amu

It is the first column that contains the *name* that we use globally within VAMDC for a certain bit of information. This is what we mean in the following when we talk about “global names” or “keywords”.

The full VAMDC dictionary is still being worked on and it currently resides at <http://vamdc.tmy.se/dict/> where also some helper tools are provided.

At the nodes, the dictionary is used in the following different ways. Note that some keywords do not make sense being used in all three cases. Common sense applies.

Note: The Returnables and Restrictables, as described in the following, are different for each node (depending on the data it offers and its structure) and need to be written when setting up a new node.

3.3.1 Returnables

Each node keeps a list of global names that we call the *Returnables*. This list contains the names associated with the kinds of information that a node has to offer. This list is offered as XML at the *tap/capabilities/* URL end point which allows user applications to decide whether it is worth to query a certain node for a certain bit of data, or not.

The node software stores the Returnables not only as a list of global names, but as a list of key-value pairs where the names are the keys and the values are the corresponding places of the data in the *data model* (see above). This way, the Returnables become a simple one-to-one map between the global names, used by all VAMDC nodes, and the node-specific layout of the database.

This “translation” is then used, among other things, by the code that fills the data into a certain output format which in turn can become node-independent. Thereby each Returnable corresponds to a certain place (a column in table format, or a certain XML tag) in the output format.

3.3.2 Restrictables

It is the list of global names that make sense to put constraints on at a certain node and therefore tells which names from the dictionary can be used in the WHERE-clause of a query to the node (see query language below).

Again, the node software uses the Restrictables as a list of key-value pairs where the keys are the global names and the values are the corresponding place in the data model. As for the Returnables, this one-to-one map of global names to custom data model allows to translate between the two - this time when the query is parsed at arrival. The code for parsing the query uses this and can thus be re-used by all nodes without altering the code.

3.3.3 Requestables

Requestables are a third way of using the dictionary. They are used in the SELECT-clause of the SQL expression when one wants to receive only a subset of the data that matches the restrictions. For example, *SELECT Species, RadiativeTransitions* would return only the fields in this group and skip any information about the states, if it were available.

Note: This is used for a future feature of the query language that is not yet implemented in the node software.

3.4 The registry

The registry is a central web service where all VAMDC nodes are registered with their access URL and some additional information. This allows finding nodes before sending queries to them. You will need to register your node there once the setup is complete.

Note: What follows below is not necessary to know for setting up a new VAMDC node.

3.5 TAP services

TAP stands for *Table Access Protocol* and is a Virtual Observatory standard definition of a web service. The detailed specs can be found [here](#). All VAMDC nodes offer their data through a TAP-like interface which means that the URL end-points are named like in TAP, the most important being */tap/sync* for a data query which returns the data synchronously (in the immediate reply). Also the attribute names for submitting a query are strongly inspired by TAP so that a query to a single VAMDC node looks something like this:

```
http://domain.of.your.node/tap/sync/?LANG=VSS1&FORMAT=XSAMS&QUERY=query string
```

VAMDC nodes currently only use and support a subset of the TAP standard, i.e. that parts that are needed within the VAMDC. Keep in mind that users will not primarily query an individual node but use a higher level tool like the VAMDC portal for querying many nodes at once. Data providers that want to set up their own VAMDC node do not really need to care about TAP either.

The more detailed specification of the VAMDC variant of a TAP service can be found in the standards documentation at <http://vamdc.org/documents/standards/>.

3.6 The query language

The node software uses the *VAMDC SQL-subset 1 (VSS1)* and will implement the future iterations of the VAMDC query language. VSS1 is basically a SQL-like string where the layout of the database behind the answering node does not need to be known - instead one uses the keywords from the dictionary in the WHERE part to restrict the selection of data. This means that all nodes understand identical queries and there is no need to adapt the query to a certain node.

Details can be found in the VAMDC-TAP specification (see link above) and should not be necessary to know for setting up a new VAMDC node. Defining the Restrictables and Returnables is enough for allowing the node software to take care of the rest.

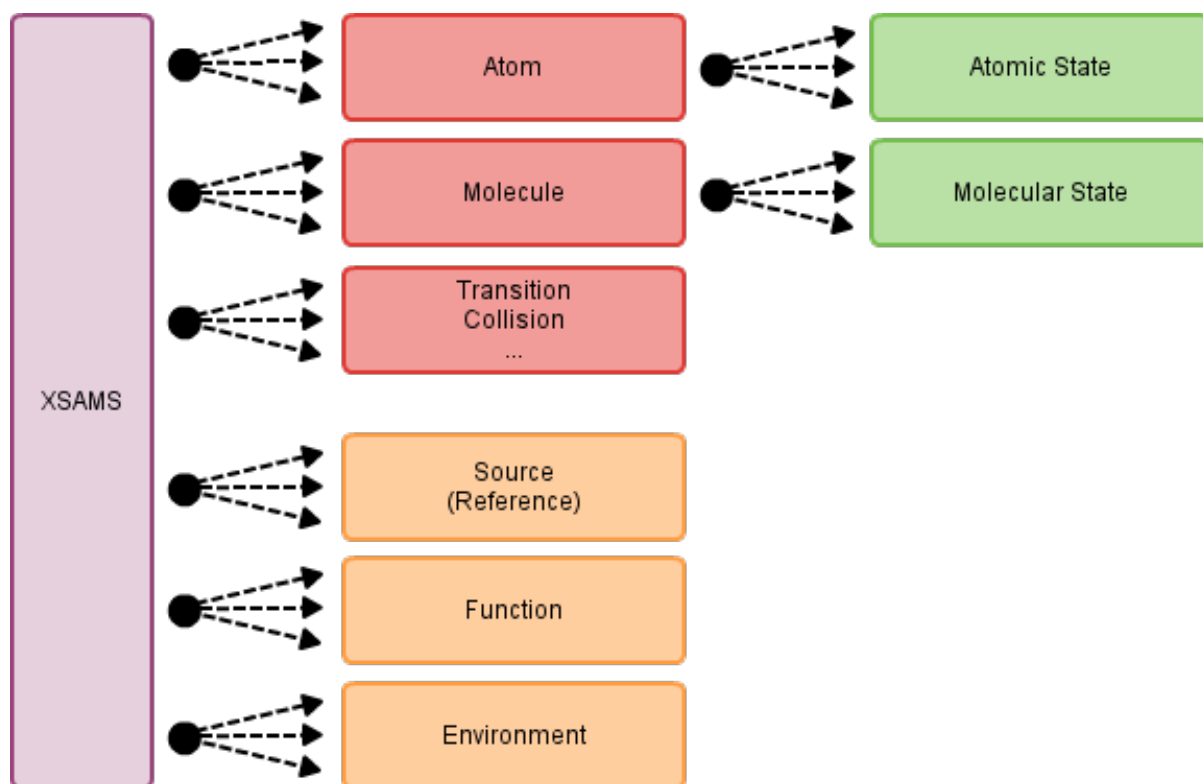
3.7 The XSAMS schema

XSAMS stands for XML Schema for Atoms, Molecules and Solids. It defines a strict way to represent data in XML. XSAMS is the format in which VAMDC nodes send their data replies.

Link to the [VAMDC-XSAMS project on Sourceforge](#).

The NodeSoftware provides an implementation of the XSAMS schema and data providers need not know it in detail to set up a VAMDC node. However, basic knowledge of its structure is needed to be able to write the few bits of code as explained in the next chapter.

XSAMS is a hierarchical structure which simplified looks like this:



Inside each box resides all the data that corresponds to it. The *atom* box holds the name, atomic number, masses, isotopes, ionization and so on. The *atomic state* box holds the state energy, quantum numbers and so on.

The different parts are interlinked, for example each atomic state has an ID and the transitions can refer to them as their initial and final state. Sources (i.e. publications) can be referenced, again by their ID, for each bit of information provided.

3.8 The generic XSAMS generator

The node software comes with an implementation of the XSAMS that can (but need not necessarily) be used by all nodes, aka the XSAMS *generator*. This frees data providers from the need to know about XML and the details of the schema. In order for this to work, data providers need fill the *Returnables* as described above and in the next chapter. The generator then knows how to put the data into the schema.

In principle, XSAMS allows many more nested loops than are shown in the diagram above. But since each node needs to build from its database the structure that matches the hierarchy, we have made some deliberate simplifications. For example we treat each ion and isotope or an atom as a different atom/species. This means that skip the complexity of having five or more nested loops at the expense of replicating some information.

3.9 The portal

The portal is the obvious example of a *user application* that makes use of VAMDC nodes. It is a web site that facilitates the submission of a query to many nodes at once by providing a web form out of which it assembles the query string which it then sends to one or many nodes, gathers the results from each of them and presents them to the user.

SOFTWARE PREREQUISITS AND INSTALLATION

4.1 Quick start

If you use a Linux-distribution like Debian (squeeze) or Ubuntu (some not too old version), you can simply run the following command (with root-rights) to install all software that you need:

```
$ apt-get update && apt-get install python python-pip python-pyparsing python-mysqldb gunicorn ng
$ pip install django
```

This will automatically install some more packages that the above ones depend upon. There are most probably similar packages for other linux distributions than Debian. All software should be able to be installed on Windows and OSX as well but it probably involves some more effort and we unfortunately cannot give support for this.

We also provide a virtual machine appliance with Debian/Linux and all required software installed into it. You can then run this virtual machine on a host computer, using VirtualBox which is available for free on most operating systems. See *Obtain and use a Virtual Machine with the NodeSoftware* for more detail on this.

If the commands above worked or you run the virtual machine, you might want to skip to *Test your installation*. Otherwise continue reading for a list of the individual software dependencies..

4.2 Python plus some modules

Python is a wide-spread, open-source, object-oriented, dynamically-typed, interpreted programming language. You can read all about it at <http://python.org> and there exist installation packages for all operating systems and architectures.

We require Python between (and including) versions 2.5 and 2.7.

We recommend to also install IPython (<http://ipython.scipy.org/>), an improved interactive shell for Python.

4.2.1 Database access library

This depends on your choice of database engine (see below). The two choices we support primarily are SQLite (access library comes with Python itself) and MySQL (access library at <http://pypi.python.org/pypi/MySQL-python/> but preferably installed by your OS's package manager).

4.2.2 PyParsing

This is needed for our SQL-parser and you can read about it at <http://pypi.python.org/pypi/pyparsing>

Again, it is best installed via your distribution's package manager.

4.3 Django

Django is the Python-based web-framework that we use to run the services (see *Introduction* and <http://djangoproject.com>). We currently use Django 1.3.X (where X is the latest bug-fix version number) but newer versions will be supported as they are released.

The packaged version of your OS is probably outdated. This is why we recommend to install Django using `pip` (see command above). Alternatively follow the installation instructions on the Django website.

4.4 Database engine

If the data that your node should serve reside already in a relational database, there is most probably no need to set up a new one but you instead deploy the node software directly on top of the existing database. The list of databases that Django can handle can be found at <http://docs.djangoproject.com/en/1.3/ref/databases/>

When setting up a new database, we recommend one of the following two

- SQLite <http://www.sqlite.org/>
- MySQL <http://mysql.com/> (or, if ORACLE succeeds in messing MySQL up, the MySQL fork called MariaDB <http://mariadb.org/>)

Unless the data set is extremely large and/or complex, the choice between the two is of minor importance. SQLite has the advantage of not relying on a separate server software and is often on par with MySQL in terms of speed. Its limitation in terms of concurrent write access is not relevant in our typical use case where the database is only read, not written to, during standard operation.

4.5 Webserver

The node software needs to run within a webserver. The two setups that we successfully tested are *Gunicorn* (together with *nginx*) and the Apache webserver (with its WSGI module).

This is covered in more detail in *Deployment of your node*.

4.6 Git version control

This is not a real requirement since you can download the node software (see *The Code*) directly. However, using the version control system *git* (<http://git-scm.com/>), it becomes easier to update your installation and to re-submit your changes.

4.7 The node software itself

See *The Code* on how to obtain the source code.

4.8 Test your installation

None of the following commands should give you an error:

```
$ python -c "import django"
$ python -c "import pyparsing"

$ cd /path/to/where/you/downloaded/NodeSoftware
$ cd nodes/ExampleNode
```

```
$ ./manage.py
$ ./manage.py test
$ ./manage.py shell
```

The last command will open an interactive Python shell for you (IPython, if you have it installed, otherwise standard Python) and in there you should be able to run:

```
>>> from node.models import *
>>> import vamdctap
>>> exit()
```

If any of this fails, please make sure you have installed all of the above correctly and ask your system administrator for help. For contacting us, see *Bugs and Contact*.

Note: The above only tests that you have installed the software correctly, not the setup and configuration of the node in question.

UPGRADING

5.1 NodeSoftware

The simplest way is to simply download the latest tar.gz-archive and extract it on top of your previous installation. We however strongly recommend to backup the files in your node-directory before doing this; alternatively moving the old NodeSoftware to a different location and then copy the files you need from there into the new version.

If you instead use our version control system, please see *Collaborating with git and GitHub* on how to get the latest.

Note: After upgrading the NodeSoftware, you should check that your node is still running properly. We cannot (yet) guarantee that you need not update your node-specific code to fit the latest version. Larger changes will be mentioned in the *Changelog*.

5.2 Django

This depends on how you installed Django. With `pip` it is enough to run:

```
$ pip install --upgrade django
```

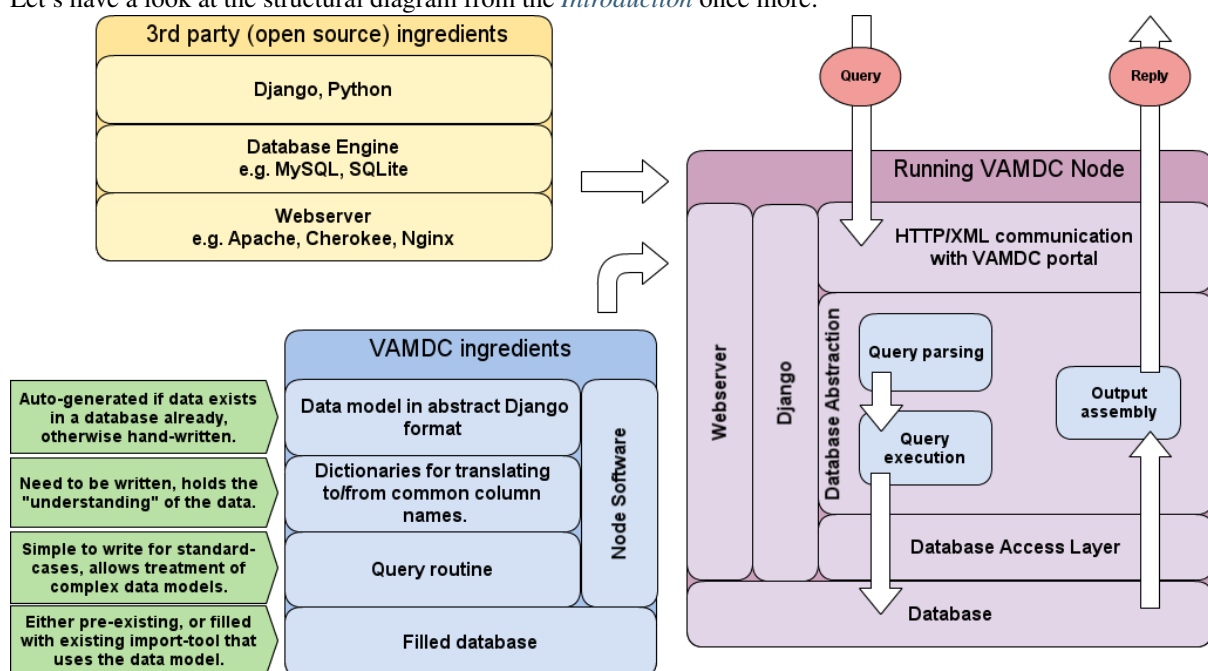
5.3 Everything else

If you have installed all the prerequisites from Debian or Ubuntu packages as recommended, you can simply run the following regularly to keep your system up to date:

```
$ apt-get update  
$ apt-get upgrade
```

STEP BY STEP GUIDE TO A NEW VAMDC NODE

Let's have a look at the structural diagram from the *Introduction* once more:



If you have followed the instructions of the page on *Software prerequisites and installation*, you are done with the yellow box in the figure. This page will tell you first how to configure and write the few code bits that your node needs before running (blue box), and then how to deploy the node and make it run as shown in the violet box.

It goes like this:

- Get the Nodsoftware and make a copy of the example node.
- Auto-create a new settings file and put your database connection there.
- **Either**
 - Write your data model and let Django create the database from it. Then use the import tool to put your data there.
 - Let Django write the model from an existing database that you already have.
- Assign names from the VAMDC dictionary to your data to make them globally understandable.
- Start your node and test it.

But let's take it step by step:

6.1 The main directory of your node

Let's give the directory which holds your copy of the NodeSoftware a name and call it `$VAMDCROOT`. (It is called *NodeSoftware* by default and exists wherever you downloaded and extracted it, unless you moved it elsewhere and/or renamed it, which is no problem to do) a name and call it `$VAMDCROOT`. Let's also assume the name of the dataset is *YourDBname*.

Inside `$VAMDCROOT` you find several subdirectories. For setting up a new node, you only need to care about the one called `nodes/` which contains the files for several nodes already, plus the example node. The first thing to do, is to make a copy of the `ExampleNode`:

```
$ export VAMDCROOT=/your/path/to/NodeSoftware/  
$ # (the last line is for Bash-like shells, for C-Shell use `setenv` instead of `export`  
$ cd $VAMDCROOT/nodes/  
$ cp -a ExampleNode YourDBname  
$ cd YourDBname/
```

Note: In the following you always work within this newly created directory for your node. You should not need to touch any files or run commands outside it.

6.2 Inside your node directory

The first thing to do inside your node directory is to run:

```
$ ./manage.py
```

This will generate a new file `settings.py` for you. This file is where you override the default settings which reside in `settings_default.py` (which you should not edit!). There are only a few configuration items that you need to fill

- The information on how to connect to your database.
- A name and email address for the node administrator(s).
- Example queries that makes sense with your data.
- Optionally you can set the location of the log-file and override other options by copying from `settings_default.py`.

The structure for filling in this information is already inside the newly created file. You can leave the default values for now, if you do not yet know what to fill in.

There are only three more files that you will need to care about in the following:

- `node/models.py` is where you put the data model,
- `node/dictionaries.py` is where you put the dictionaries and
- `node/queryfunc.py` is where you write the query function,

all of which will be explained in detail in the following.

6.3 The data model and the database

By *data model* we mean the piece of Python code that tells Django the layout of the database, including the relations between the tables. By *database* we mean the actual relational database that is to hold the data. (See also *The main concepts behind the implementation*).

There are two basic scenarios to come up with these two ingredients. Either the data are already in a relational database, or you want to create one.

6.3.1 Case 1: Existing database

If you want to deploy the VAMDC node software on top of an existing relational database, the *data model* for Django can be automatically generated by running:

```
$ ./manage.py inspectdb > node/models.py
```

This will look into the database that you told Django about in `settings.py` above and create a Python class for each table in the database and attributes for these that correspond to the table columns. An example may look like this:

```
from django.db.models import *

class Species(Model):
    id = IntegerField(primary_key=True)
    name = CharField(max_length=30)
    ion = IntegerField()
    mass = DecimalField(max_digits=7, decimal_places=2)
    class Meta:
        db_table = u'species'
```

There is one important thing to do with these model definitions, apart from checking that the columns were detected correctly: The columns that act as a pointer to another table need to be replaced by *ForeignKeys*, thereby telling the framework how the tables relate to each other. This is best illustrated in an example. Suppose you have a second model, in addition to the one above, that was auto-detected as follows:

```
class State(Model):
    id = IntegerField(primary_key=True)
    species = IntegerField()
    energy = DecimalField(max_digits=17, decimal_places=4)
    ...
```

Now suppose you know that the field called *species* is actually a reference to the *species*-table. You would then change the class *State* as such:

```
class State(Model):
    id = IntegerField(primary_key=True)
    species = ForeignKey(Species)
    energy = DecimalField(max_digits=17, decimal_places=4)
    ...
```

Note: You will probably have to re-order the classes inside the file `models.py`. The class that is referred to needs to be defined before the one that refers to it. In the example, *Species* must be above *State*.

Let's add a third model:

```
class Transition(Model):
    id = IntegerField(primary_key=True)
    species = ForeignKey(Species)
    upper_state = ForeignKey(State, related_name='transup')
    lower_state = ForeignKey(State, related_name='translo')
    wavelength = FloatField()
```

The important thing here is the *related_name*. Whenever you want to define more than one *ForeignKey* to the *same* model, you need to set this to an arbitrary name. This is because Django will automatically set up the reverse key for you and needs to give it a unique name. The reverse key in this example could be used to get all the *Transitions* that have a given *State* as upper or lower state. More on this at [Setting the related name of a field](#).

Once you have finished your model, you should test it. Continuing the example above you could do something like:

```

$ ./manage.py shell
>>> from node.models import *
>>> allspecies = Species.objects.all()
>>> allspecies.count() # the number of species is returned
>>> somestates = State.objects.filter(species__name='He')
>>> for state in somestates: print state.energy
>>> sometransitions = Transition.objects.filter(wavelength__lt=500)
>>> atransition = sometransitions[5]
>>> othertransitions = atransition.upper_state.transup.objects.all()
>>> othertransitions.count() # gives the number of transitions with the
                             # same upper state.

```

Detailed information on how to use your models to run queries can be found in Django's own excellent documentation: <http://docs.djangoproject.com/en/1.3/topics/db/queries/>

6.3.2 Case 2: Create a new database

In this case we assume that the data are in ascii tables of arbitrary layout. The steps now are as follows:

1. Write the data model in your `node/models.py`.
2. Create an empty database with corresponding user and password
3. Tell the node software where to find this database.
4. Let the node software create the tables
5. Use the import tool to fill the database with the data.

First of all, you need to think about how the data should be structured. Data conversion (units, structure etc) can and should be done while importing the data since this saves work and execution time later. Since the data will need to be represented in the common XSAMS format, it is recommended to adopt a layout with separate tables for species, states, processes (radiative, collisions etc) and references.

Deviating data models are certainly possible, but will involve some more work on the query function (see below). In any case, do not so much think about how your data is structured now, but how you want it to be structured in the database, when writing the models.

Writing your data models is best learned from example. Have a look at the example from Case 1 above and at file `$VAMDCROOT/nodes/vald/node/models.py` inside the NodeSoftware to see how the model for VALD looks like. Keep in mind the following points:

- As mentioned, a *class* in the model becomes a *table* in the database and the fields/members of the class correspond to the table columns.
- Each class should have one member with `primary_key=True`. If not, one called *id* will be implicitly created for you.
- How you name your classes and fields is up to you. Sensible names will make it easier to write the dictionaries below.
- Use the appropriate field type for each bit of data, e.g. `BooleanField`, `CharField`, `PositiveSmallIntegerField`, `FloatField`. There is also a `DecimalField` that allows you to specify arbitrary precision which will also be used in later ascii-representations of data.
- Use `ForeignKey()` to another class's primary key to connect your tables.
- The full list of possible fields can be found at <http://docs.djangoproject.com/en/1.3/ref/models/fields/>.
- If you know that a field will be empty sometimes, add `null=True` to the field definition inside the brackets `()`.
- For fields that are frequent selection criteria (like wavelength for a transition database), you can add `db_index=True` to the field to speed up searches along this column (at the expense of some disk space and computation time at database creation).

- If you do not define a table name for your model with the Meta class, as in the first example above, the table in the database will be named as the model, but lowercase and with a prefix *node_*.

Once you have a first draft of your data model, you test it by running (inside your node directory):

```
$ ./manage.py sqlall node
```

This will (if you have no error in the models) print the SQL statements that Django will use to create the database, using the connection information in `settings.py`. If you do not know SQL, you can ignore the output and move straight on to creating the database:

```
$ ./manage.py syncdb
```

Now you have a fresh empty database. You can test it with the same commands as mentioned at the end of Case 1 above, replacing “Species” and “State” by your own model names.

Note: There is no harm in deleting the database and re-creating it after improving your models. After all, the database is still empty at this stage and *syncdb* will always create it for you from the models, even if you change your database engine in `settings.py`. The command for re-creating the tables in the database (deleting all data!) is `./manage.py reset node`.

Note: If you use MySQL as your database engine, we recommend its internal storage engine InnoDB over the standard MyISAM. You can set this in your `settings.py` by adding `'OPTIONS': {"init_command": "SET storage_engine=INNODB"}` to your database setup. We also recommend to use UTF8 as default in your MySQL configuration or create your database with `CREATE DATABASE <dbname> CHARACTER SET utf8;`

How you fill your database with information from ascii-files is explained in the next chapter: [How to get your data into the database](#). You can do this now and return here later, or continue with the steps below first.

6.4 Using the XML generator

Before we go on to the remaining two ingredients, the *query function* and the *dictionaries*, we need to have an understanding on how they play together in the XML generator. As you remember from [The XSAMS schema](#), the goal is to run queries on your models and pass on the output to the generator so that it can looped over them to fill the hierarchical XSAMS structure.

In order to make this work, we need to name the variables that you pass into the generator (as explained below) and the loop variables that you use in the Returnables. For example, continuing on the model above: Assume you have made a selection of your Transition model; you pass this on under the name *RadTrans*; the generator loops over it, calling each Transition inside its loop *RadTran* (note the singular!). *RadTran* is now a single instance of your Transition model and has the wavelength as *RadTran.wavelength* since we called the field this way above. The entry in the RETURNABLES would therefore look like `'RadTranWavelength': 'RadTran.wavelength'` - where the first part is the keyword from the VAMDC dictionary (which the generator knows where in the schema it should end up) and the second part tells it how to get the value from the query results that it got from your query function.

Do not fret if this sounded complicated, it will become clear in the examples below. Just read the previous paragraph again after that.

Here is a table that lists the variables names that you can pass into the generator and the loop variables that you use in the Returnables. The one is simply the plural of the other.

Passed into generator	Loop variable	Object looped over	Loop variable
Atoms	Atom	Atom.States Atom.SuperShells Atom.Shells	AtomState AtomSuperShell AtomShell
Continued on next page			

Table 6.1 – continued from previous page

Molecules	Molecule	Atom.ShellPairs	AtomShellPair
		Molecule.States	MoleculeState
Solids	Solid	Solid.Layers	Layer
		Solid.Layers.Components	Component
Particles	Particle		
RadTran	RadTran	RadTran.ShiftingParams	ShiftingParam
		RadTran.ShiftingParams.Fits	Fit
		RadTran.ShiftingParams.Fits.Parameters	Parameter
RadCross	RadCros		
		RadCros.BandModes	BandMode
CollTran	CollTran		
		CollTran.Reactants	Reactant
		CollTran.IntermediateStates	IntermediateState
		CollTran.Products	Product
		CollTran.DataSet	DataSet
		CollTran.DataSet.FitData	FitData
		CollTran.DataSet.FitData.Arguments	Argument
		CollTran.DataSet.FitData.Parameters	Parameter
		CollTran.DataSet.TabData	TabData
NonRadTran	NonRadTran		
Environments	Environment		
		Environment.Species	EnvSpecies
Particles	Particle		
Sources	Source		
Methods	Method		
Functions	Function		
		Function.Parameters	Parameter

The third and fourth columns are for an inner loop. So for example the generator loops over all *Atoms*, calling each atom instance *Atom*. To extract all states being a part of this particular *Atom*, the generator will assume that there is an iterable *States* defined on each *Atom* over which it will iterate. So it will loop over *Atom.States*, calling each of state *AtomState* in the inner loop, like this:

```
for Atom in Atoms:
    [...]
    for AtomState in Atom.States:
        [...]
```

It is up to you to make sure the *Atom.States* is defined if you want to output state information. This is covered in the next section.

6.5 The query routine

Now that we have a working database and data model and know in principle how the generator works, we simply need to tell the framework how to run a query and pass the output to the generator. This is done in a single function called *setupResults()* which must be written in the file `node/queryfunc.py` in your node directory. It works like this:

- *setupResults()* is called from elsewhere and you need not run it yourself.

- `setupResults()` gets an object as input, called `sql`. This is a parsed version of the query that comes in. It holds the WHERE-part as `sql.where` and so on.
- We now need to run this query on the data model in order to get so called *QuerySets* which are basically unevaluated queries that are then passed on to the XML generator which takes care of the rest.
- If you want to enforce limits on how much data can be returned in one query, this can be done here as well.
- You should also calculate some statistics on how much information a query returns and return it as header information.

In a concrete example of an atomic transition database, it looks like this:

```

from django.db.models import Q
from vamdc.tap.sqlparse import *
from dictionaries import *
from models import *

LIMIT = 10000

def setupResults(sql):
    q = eval( where2q(sql.where, RESTRICTABLES) )
    transs = Transition.objects.filter(q).order_by('wavelength')
    ntranss = transs.count()

    if ntranss > LIMIT:
        percentage = '%.1f'%(float(LIMIT)/ntranss *100)
        limitwave = transs[LIMIT].wavelength
        transs = Transition.objects.filter(q,Q(vacwave__lt=limitwave))
    else: percentage=None

    spids = set( transs.values_list('species_id', flat=True) )
    species = Species.objects.filter(id__in=spids)
    nspecies = species.count()
    nstates = 0
    for specie in species:
        subtranss = transs.filter(species=specie)
        up=subtranss.values_list('upper_state_id', flat=True)
        lo=subtranss.values_list('lower_state_id', flat=True)
        sids = set(up+lo)
        specie.States = State.objects.filter(id__in = sids)
        nstates += len(sids)

    headerinfo={'TRUNCATED':percentage,
                'COUNT-ATOMS':nspecies,
                'COUNT-STATES':nstates,
                'COUNT-RADIATIVE':ntranss
                'APPROX-SIZE':ntranss*0.001
               }

    return {'RadTrans':transs,
            'Atoms':species,
            'HeaderInfo':headerinfo
           }

```

Explanations on what happens here:

- Lines 1-4: We import some helper functions from the sqlparser and the dictionaries and models that reside in the same directory as `queryfunc.py`
- Line 6: Set the limit of transitions for use below.
- Line 7: Begin the function `setupResults`. Do not change this line.
- Line 9: This uses the helper function `where2q()` to convert the information in `sql.where` to `QueryObjects` that match your model, using the `RESTRICTABLES` (see below). The result from `where2q()` is a string that

needs to be executed with `eval()`.

- In line 10 we simply pass these `QueryObjects` to the Transition model's filter function. This returns a `QuerySet`, an unevaluated version of the query, which we assign to the variable `transs`. We also ordered it by `wavelength`.
- Line 11: We use the `count()` method on the `QuerySet` to get the number of transitions which we later pass into the header.
- Line 13-17: We check if the number is larger than our limit and shorten the `QuerySet` if necessary. This is done by getting the `wavelength` at the limit and making a new `QuerySet` that has as an additional restriction the new upper `wavelength` limit. We also prepare a string with the percentage for the headers.
- Lines 19-29: Here comes the tricky part. For the selected transitions, we now need to create the corresponding atoms/species, since they go into different parts of the generator, see the table above. Not only that, each atom should have attached its list of states that are upper or lower states for the selected transitions - there is an inner loop over `Atom.States` in the generator, remember? In detail:
 - Line 19: We pull a single column out of the Transitions model, the key that links to the Species model. We put that into a `set()` to throw out duplicates.
 - Line 20: We use this set to query for all our Species.
 - Line 21: We count them and save the result for later.
 - Line 22: We make a new variable for the number of states which we will increase in the coming loop.
 - Line 23: Start a loop over our selected `species`.
 - Line 24: Make a sub-selection on our previously selected transitions, now only selecting the ones that belong to the current species.
 - Lines 25-26: As for the species IDs before, we now pull the keys to the upper and lower states out of our Transition model.
 - Line 27: We concatenate the two lists of IDs and put them in a `set()` to get rid of duplicates. `sids` is now a list of IDs of all the states within the current species that are used in the selected transitions.
 - Line 28: Use this list to make the query on the State model. And, **most importantly**, attach it to the current species object. This way we have constructed the nested structure for the generator.
 - Line 29: For the statistics, we now increase the state count with the number for the current species.
- Lines 31-36: Put the statistics into a key-value structure where the keys are the header names as defined by the VAMDC-TAP standard and the values are the strings/numbers that we calculated above.
- Lines 39-41: Return the `QuerySets` and the headers, again as key-value pairs. The keys are the names from the first column of the table above, so that the generator recognizes them and loops over them at the right place.

Note: As you might have noticed, all restrictions are passed to the Transitions model in the above example. This does not mean that we cannot put constraints on e.g. the species here. We simply use the models `ForeignKey` in that case in the `RESTRICTABLES`. An entry there could e.g. be `'AtomIonCharge': 'species__ion'` which will use the `ion` field of the species model. Depending on your database layout, it might not be possible to pass all restrictions to a single model. Then you need to write a more advanced query than the shortcuts in Lines 7-8.

Note: We are well aware that adapting the above example to your data is a non-trivial task unless you know Python and Django reasonably well. There is a more complete example in `ExampleNode/node/queryfunc.py` and you can also have a look at the other nodes' `queryfunc.py` which are included in the NodeSoftware. And, of course, we are willing to assist you in this step, so feel free to contact us about this.

More comprehensive information on how to run queries within Django can be found at <http://docs.djangoproject.com/en/1.3/topics/db/queries/>.

6.6 The dictionaries

As the last important step before the new node works, we need to define how the data relates to the VAMDC *dictionary*. If you have not done so yet, please read *The VAMDC dictionary* before continuing.

What needs to be put into the file `node/dictionaries.py` is the definition of two variables that map the individual fields of the data model to the names from the dictionary, like this:

```
RESTRICTABLES = { \
'AtomSymbol' : 'species__name' ,
'AtomIonCharge' : 'species__ion' ,
'RadTransWavelength' : 'wavelength' ,
}

RETURNABLES={ \
'NodeID' : 'YourNodeName' , # constant strings work
'AtomIoncharge' : 'Atom.ion' ,
'AtomSymbol' : 'Atom.name' ,
'AtomStateEnergy' : 'AtomState.energy' ,
'RadTransWavelength' : 'RadTran.wavelength' ,
}
```

Note: Note for example the use of the names *Atom* and *AtomState* on the right-hand side of the dictionary definition. These are examples of the “loop variables” mentioned in the table above and act as shortcuts to the nested data you are storing.

There are tools for getting started with writing these and for validation once you are done at <http://vamdc.tmy.se/dict/>

6.6.1 About the RESTRICTABLES

As we have learned from writing the query function above, we can use the RESTRICTABLES to match the VAMDC dictionary names to places in our data model. The key in each key-value-pair is a name from the VAMDC dictionary and the values are the field names of the model class that you want to query primarily (Transition, in the example above, line 10).

The RESTRICTABLES example give fits our query function from above, so we know that the “main” model is the Transitions. Now if a query like “AtomIonCharge > 1” comes along, this can be translated into *Transition.objects.filter(species__ion__gt=1)* without further ado, which is exactly what *where2q()* does. Note that we here used a ForeignKey to the Species model; the values in the RESTRICTABLES need to be written from the perspective of the queried model.

Note: Even if you chose to not use the RESTRICTABLES in your `setupResults()` and treat the incoming queries manually, you are still encouraged to fill the keys (with the values being empty), because they are automatically provided to the VAMDC registry so that external services can figure out which names make sense to query at this node.

6.6.2 About the RETURNABLES

Equivalent to how the RESTRICTABLES take care of translating from global names to your custom data model when the query comes in, the RETURNABLES do the opposite on the way back, i.e. when the data reply is sent by the generator, as we have already seen above.

Again the keys of the key-value-pairs are the global names from the VAMDC dictionary. The values now are their corresponding places in the QuerySets that are constructed in `setupResults()` above. This means that the XML generator will loop over the QuerySet, getting each element, and try to evaluate the expression that you put in the RETURNABLES.

Continuing our example from above, where the State model has a field called *energy*, so each object in the Query-Set will have that value accessible at *AtomState.energy*. Note that the first part before the dot is not the name of your model, but the *loop variable* inside the generator as it is listed in the second (or forth, in the case of an inner loop) column of the table above.

There is only one keyword that you must fill, all the others depend on your data. The obligatory one is *NodeID* which you should set to a short string that is unique to your node. It will be used in the internal reference keys of an XSAMS document. By including the NodeID, we make these keys globally unique within VAMDC which will facilitate the merging of data that come from different nodes.

It was mentioned before, but now is the time to point you once more to <http://vamdc.tmy.se/dict/> where you first of all can browse all the available keywords. By selecting the ones that match your data, you can download a raw version of your `dictionaries.py` which you then fill in. The website also can perform some tests on your Returnables and Restrictables for finding errors.

Note: Again, at least the keys of the RETURNABLES should be filled (even if you use your own generator for the XML output) because this allows the registry to know what kind of data your node holds before querying it.

6.7 Testing the node

Now you should have everything in place to run your node. If you still need to fill your database with the import tool, now is the time to do so according to *How to get your data into the database*.

Django comes with a built-in server for testing. You can start it with:

```
$ ./manage.py runserver
```

This will use port 8000 at your local machine which means that you should be able to browse to <http://127.0.0.1:8000/tap/availability/> and hopefully see a positive status message.

You should also be able to run queries by accessing URLs like:

```
http://127.0.0.1:8000/tap/sync?LANG=VSS1&FORMAT=XSAMS&QUERY=SELECT ALL WHERE AtomIonCharge > 1
```

replacing the last part by whatever restriction makes sense for your data set.

Note: The URL has to be URL-encoded when testing from a script or similar. Web browsers usually do that for you. To also see the statistics headers, you can use `wget -S -O output.xml "<URL>"`.

A more extensive test framework is in the making and will be documented here soon. In any case you should run test queries to your node and make sure that the output in terms of volume and values matches your expectations.

Once your node does what it should with the test server, you can start thinking about *deploying it*.

HOW TO GET YOUR DATA INTO THE DATABASE

In the previous chapter, we have learned how to define the database layout and tell the framework to create the database accordingly. The following describes how to fill this database with data that reside one or many ascii tables.

Note: There are many ways to achieve this and you are certainly free to fill the database in any way you want, if you already know how to do it.

The strategy we adopt is to use the database's own import mechanisms which are manyfold faster for large amounts of data than manually inserting row by row.

So the import becomes a two-step process:

1. create one ascii file per data model where each of them has columns that exactly match the columns in the database.
2. run one SQL command for each of these files to load it into the matching database table.

Since you might already have step 1 finished or might be able to get it with your own data handling tools, let's have a look at step 2 first.

7.1 Loading ascii data into the database

In the following, we assume that you use MySQL as your database engine which is our recommendation when a new database is set up for the first time. Other engines have similar mechanisms for bulk loading data.

The command we use looks like this:

```
mysql> LOAD DATA INFILE '/path/to/data.file' into table <TAB>;
```

where <TAB> is the name of the table that matches the file that is loaded.

Note: The table names have a prefix *node_*, i.e. the table for a model called *State* will be called *node_state*, unless you specify the table name in the model's definition. You can see a list of all tables by running *SHOW TABLES*;

The LOAD DATA command has several more options and switches for setting the column delimiter, skipping header lines and the like. Mathematical or logical operations can be run on the columns, too, before the data get inserted into the database.

You can read all about LOAD DATA at <http://dev.mysql.com/doc/refman/5.1/en/load-data.html>

A more complete example would look like:

```
mysql> LOAD DATA INFILE '/path/transitions.dat' IGNORE INTO TABLE transitions COLUMNS TERMINATED
```

7.2 Preparing the input files

In the not so unlikely case that the data are not yet in a format that exactly matches the database layout, we provide a tool that can be used to re-write your data and create the ascii files that can be loaded as described above.

These files must fulfill the following criteria:

- One file per database table. LOAD DATA cannot update existing rows.
- Same number of columns in the file as in the table and in the right order. Although LOAD DATA can take a list of columns to circumvent this restriction, it makes sense to get this right.
- Links between the tables are in place. The key values that link tables (e.g. states and transitions) should be already in the ascii files (even though they *can* still be generated with LOAD DATA by using some SQL magic).
- A consistent delimiter between the columns (no fixed record length) and consistent quoting.
- Empty (NULL) values are written as W, not 0 or anything else. (Can also be fixed later if this is the only thing missing)

The Node Software ships with a *rewrite tool* for creating the files according to these criteria. The tool can be used to import almost any format of file as long as it stores its records as *lines* (blocks of data stretching several lines in the file are currently only partly supported).

For using the rewrite tool, you need to tell it how your original data files are named and how they are structured. This is done in something called a *mapping file*. The mapping file describes how the rewriter should extract data from your custom text files and write them into the files that match the data model.

should usually import helper functions from *imptools/linefuncs.py* to do much of the work for you. There is a sample mapping file in the *imptools/* directory, you can copy that to your Node to edit.

Starting the rewrite

Once you have defined the mapping file as described in the following section, you give it as an argument to the *imptools/run_rewrite.py* program:

```
$ python run_rewrite.py mapping_mynode.py
```

7.3 The mapping file

The mapping file is a standard Python file. It must define a variable called *mapping* which contains a list of definitions that describe how the rewriter should parse the lines of any number of text files and put the result into the output files.

Let's start by defining your input files:

```
from imptools.linefuncs import *

# the names of the input files
basepath = "/path/to/your/raw_data/"
file1 = basepath + 'raw_file1.txt'
file2 = basepath + 'raw_file2.txt'
file3 = basepath + 'raw_file3.txt'
outfile1 = basepath + 'references.dat'
outfile2 = basepath + 'species.dat'

mapping = [ ... ] # described below
```

7.3.1 The mapping list

The mapping variable is a list of Python *dictionaries*. A python dictionary is written as {key:value, key2:value2, ... }. One of these keys, *linemap*, is itself a list with further dictionaries. The structure looks like this:

```
mapping = [
  {key : value,
   key : value,
   linemap : [
     {linemap_key : value,
      linemap_key : value},
     {linemap_key : value,
      linemap_key : value}] }
  {key : value,
   key : value,
   linemap : [ ... ]}
]
```

The keys and values of each dictionary describes how to populate one output file using any number of source text files.

key	value
<i>Mandatory</i>	
outfile	The name of the file that should be created.
infile	Input file(s). If more than one file is used, this should be a list of filenames.
linemap	A list of dictionaries defining how to parse each line of the file(s) into its components.
<i>Optional</i>	
headlines	Number of header lines at the top of the input file() (default: 0).
commentchar	Which comment symbol is used in the input file(s) (default: '#').
cnull	Values in the input file(s) that should be considered 'null' and ignored (no default).
errline	Whole lines in the input file(s) that should be considered non-valid and ignored (no default).
lineoffset	An offset step length (in number of lines) between two or more read input files. Default (0) means stepping one line at a time. An offset of 1 means skipping every other line. So a lineoffset of (0,2) would mean that while every line is read in the first file, only every third is used in the second file (default is 0 offset).

If you are using more than one input file to populate one output file (for example if you read one piece of data from each file and combines them), you need to supply lists to all entries identifying features in the files, such as *commentchar*, *cnull* etc. If you do not the rewriter will return errors. Note that in order to correlate several files like this they all have to have its data in the form of lines, and be able to step systematically through those lines. Use *lineoffset* to step at different rates through the files.

The *linemap* key points to another list with dictionaries. This is the actual operating piece of code and describes exactly how to parse each line (or lines, if more than one input file is used). The result of each dictionary is the population of one database field in your model.

linemap	keyglue
<i>Mandatory</i>	
cname	The name of the field in your database model.
cbyte	A tuple (<i>linefunction</i> , <i>arguments</i>). This defines a function capable of parsing the line(s) to produce the data needed to feed to the field <i>cname</i> . The only provision of a <i>linefunction</i> is that it should take an argument <i>linedata</i> as its first argument. This contains the current line to parse, or a list of lines if more than one input files where read simultaneously.
<i>Optional</i>	
debug	This will activate verbose error messages for this parsing only. Useful for finding problems with the mapping.

Continuing our example, here's of how this could look in the mapping file (the line breaks are technically not needed, but make things easier to read):

```
mapping = [
  # first dictionary, writing into references.dat
  {
    'outfile': outfile1,
    'infiles': file1,
    'headlines' : 3,
    'commentchar' : '#',
    'linemap' : [
      {'cname': 'dbref',
       'cbyte': (bySepNr, 0, '||')},
      {'cname': 'author',
       'cbyte': (bySepNr, 1, '||')},
      # ...
    ]
  }
  # next model dictionary, writing species.dat
  {
    'outfile' : outfile2,
    'infiles' : (file2, file3), # using more than one file!
    'commentchar' : (';', '#'),
    'headliens' : (1, 3),
    'lineoffset' : (0, 1),
    'linemap' : [
      {'cname': 'pk',
       'cbyte': (charrange, 23, 25)}, # pick a range by index
      {'cname': 'mass',
       'cbyte': (charrange, 45, 45, 1)}, # retrieved from file3!
      # ...
      {'cname': 'source',
       'cbyte': (charrange, 0, 10),
      ]
    ]
  }
]
```

Here we define how to populate two models. The first dictionary makes use of the *bySepNr* line function (see below) to extract data from each line. The second instead relies on a line function called *charrange* to mix info from two input files.

7.3.2 The line functions

Since the mapping file is a normal Python module, you are free to code your own line functions to extract the data from each line in your file. There are only three requirements for how a line function may look:

- The function must take at least one argument, which holds the current line being processed, as a string. The import program will automatically send this to the function as it steps through the file. If more than one file is traversed, this input will be in the form of a *list* of line strings (it is then up to you which one to use).

- The function must return its extracted piece of data in a format suitable for the field it is to be stored in. So a function parsing data for a CharField should return strings, whereas one parsing for an IntegerField should return integer values.
- If the function is used to populate a Many-to-Many relationship (that is, the key *multireference* is set in the parsing dictionary), the line function must return a *list* of parsed results, one for each reference that is to be searched for in the database and tied to the field.

Below is a simple example of a line function that fulfills all these criteria:

```
def charrange(linedata, start, end):
    """
    Simple extractor that cuts out part of a line
    based on string index
    """
    return linedata[start:end].strip()
```

In the mapping dictionary we call this with e.g. `'cbyte' : (charrange, 12, 17)`. The first element of the tuple is the function object, everything else will be fed to the function as arguments.

This function assumes that linedata is a simple string, and so it will not work if we where to re-use it for multiple in-files (linedata will then be a list). So let's do a simple addition:

```
def charrange(linedata, start, end, filenum=0):
    """
    Simple extractor that cuts out part of line(s)
    based on string index
    """
    if is_iter(linedata):
        # this is an iterable (i.e. a list)
        # so pick one line based on linenum
        linedata = linedata[linenum]
    return linedata[start:end].strip()
```

This you can still call the same way as before, but when working with more than one file, you can also add an extra argument to pick which file to use the line from.

The import tool comes with a basic set of the most common line functions, such as extracting by line index, by separator and some more. Just `import linefuncs *` from your mapping file to make them available. You can find more info in the [Linefuncs Documentation](#).

More advanced line parsing

Sometimes you need more advanced parsing. Say for example that you need to parse two different sections of lines from one or more files and combine them into a unique identifier that you will then use as a key for connecting your model to another via a One-to-Many relationship. Or maybe you want to put a value in different fields depending on if they are bigger/smaller than a certain value. The default line functions in *linefuncs.py* cannot do this out of the box.

The solution is to write your own line function. You have the full power of Python at your command. Often you can use the default functions as “building blocks”, linking them together to get what you want. Just code your custom line functions directly in the mapping file.

Here is an example of a line function that wants to create a unique id by parsing different parts of lines from different files:

```
def get_id_from_line(linedata, seprnr, index1, index2):
    """
    extracts id from several lines.
    seprnr - nth separator to pick from file 1
    index1, index2 - indices marking piece to pick from file 2

    (file3 is always used the same way, so we hard-code the
    indices for that file.)
    """
```

```
"""
11 = bySepNr(linedata[0], sepr, ',')
12 = charrange(linedata[1], index1, index2)
13 = charrange(linedata[2], 0, 3)
if 13 == '000':
    13 = 'unknown'
# create unique id
return "%s-%s-%s" % (11, 12, 13)
```

Here we made use of the default line functions as building blocks to build a complex parsing using three different files. We also do some checking to replace data on the spot. The end result is a string combined from all sources. This would be called from the line mapping dictionary with e.g. `cbyte: (get_id_from_line, 3, 25, 29)`.

In the `imptools` directory you can find a fully functioning mapping used for importing the VALD database. It also contains a set of custom line functions to use for inspiration.

DEPLOYMENT OF YOUR NODE

Now that you have a node that runs nicely with Django's test server, the last remaining step is to configure the server that will run the node in a production setup.

How and on which server you set up your node to run permanently, is much dependent on your technical resources and the solution we give here is just one out of several possibilities (although we also quickly mention the most common alternative).

8.1 Gunicorn plus proxy

Our recommended way for hosting your node by yourself on a server is Gunicorn (<http://gunicorn.org/>, *apt-get install gunicorn* on a Debian system) which is aware of Django and understands its settings.

You would write a *gunicorn.conf* file (you find it in *nodes/ExampleNode*) like this:

```
import os
def numCPUs():
    if not hasattr(os, "sysconf"):
        raise RuntimeError("No sysconf detected.")
    return os.sysconf("SC_NPROCESSORS_ONLN")
workers = numCPUs() * 2 + 1

#bind = "127.0.0.1:8000"
bind = "unix:/tmp/gunicorn.sock"
pidfile = "/tmp/gunicorn.pid"
logfile = "/tmp/gunicorn.log"
loglevel = "info"
timeout = 60
daemon = True
```

and then simply start it from within your node directory with:

```
$ gunicorn_django -c gunicorn.conf
```

The example config makes Gunicorn listen at a unix-socket. Even though you can connect it to a TCP-port instead (see commented out line), you do not want external requests sent directly to Gunicorn, but to a proxy instead. This proxy takes care of the load balancing between the Gunicorn worker processes and can compress the XML output from your node before sending it.

8.1.1 Nginx as proxy

Nginx (<http://nginx.org/en/>, *apt-get install nginx* on a Debian system) is a fast and light-weight web server. To configure it to serve the running node with Gunicorn, according to the example above, you would configure it like this:

```

upstream app_server {
    server unix:/tmp/gunicorn.sock;
}

server {
    listen 8080; ## listen for ipv4
    listen [::]:8080 default ipv6only=on; ## listen for ipv6
    server_name your.server.domain.name;
    access_log /var/log/nginx/vamdc.access.log;

    location /yournode/tap/ {
        proxy_set_header X-Forwarded-For $proxy_add_x_forwarded_for;
        proxy_set_header X-Real-IP $remote_addr;
        proxy_set_header Host $http_host/yournode;
        proxy_pass http://app_server/tap/;
        proxy_redirect http://app_server/tap/ /yournode/tap/;

        gzip on;
        gzip_types text/plain application/xml text/xml;
        gzip_proxied any;
    }
}

```

Note that you probably want to edit the port, server name and the location at which to serve the node (change */yournode/tap* at three places but make them match each other).

If you installed *nginx* with the *debian/ubuntu* package, you can place symbolic links to the config file into */etc/nginx/* like this to make it use the config above:

```

$ cd /etc/nginx/sites-available/
$ sudo ln -s $VAMDCROOT/nodes/YourNode/nginx.conf vamdcnode
$ cd ../sites-enabled/
$ sudo ln -s ../sites-available/vamdcnode
$ sudo /etc/init.d/nginx restart

```

8.1.2 Proxy Alternatives

What you choose as proxy for Gunicorn is somewhat arbitrary. Common alternatives to *nginx* are *lighttpd* or *Apache*. Especially if the server that is to run your node already has an Apache running for serving other websites, it makes sense to simply tell it how to proxy your Gunicorn server:

```

ProxyPass /yournode http://localhost:8000
ProxyPassReverse /yournode http://localhost:8000

```

8.2 Deployment in Apache

As an alternative to deployment with Gunicorn plus proxy, the Apache webserver can not only act as a proxy but also replace Gunicorn by using its *mod_wsgi* plugin to run the Python code directly. The main disadvantage of this setup is that you cannot configure and restart the node independently from Apache, so the likelihood of interfering with any other sites that Apache serves is larger.

There are two example files in your node directory for setting this up:

- *apache.conf*: This is an Apache config file that defines a virtual server, bound to a certain host name. You will have to edit several things in that file before it will work in Apache: the server name and the path to the node software in a few places. On a Debian-like system you would then move this file to */etc/apache2/sites-available/vamdcnode* and run *a2ensite vamdcnode* to activate it.
- *django.wsgi*: This is the file that the previous one points to in its *WsgiScriptAlias*. Edit the path and your node's name.

Once you have set this up and re-started the Apache webserver, your node should deliver data at the configured URL.

8.3 Third party hosting

There are several upcoming hosting solutions that support Django directly so that you simply would upload the code and your database and everything is taken care of for you. Once these services mature, they are probably a very good solution for nodes with relatively small volumes of data.

Searching the web for “django hosting” will point you in the right direction, as does this list <https://convore.com/django-community/django-hosting-explosion/>

8.4 Logging

Finally, a few words on logging the access to your node. There are two basic ways:

- let the webserver do it.
- let the NodeSoftware do it.

The webserver/proxy, be it nginx or apache, keeps a log on when, how and by whom your node is accessed. Since the query itself is in the accessed URL, it also ends up in these logs. There are many tools to analyze and visualize this kind of logs. In the case of Apache/WSGI-deployment, errors in the NodeSoftware show up the webserver's error-log since it is the former that executes the latter. With gunicorn, the webserver knows nothing about the NodeSoftware's errors since it only acts as a proxy. Gunicorn keeps its own logs.

However, the webserver logs usually contain no information about what happened inside the NodeSoftware. If you want to keep tabs on how much data was returned from each query, how long it took to process and so on, you need to tell the NodeSoftware to save this information for you - this is where the *logging*-facility comes into play.

Nodes will primarily use this in their `queryfunc.py` where you initialize it like this:

```
import logging
>>> log = logging.getLogger('vamdc.node.queryfu')
```

Then any of the following can be used to log messages of different levels:

```
>>> log.debug('some text with a variable: %s'%variable)
>>> log.info('bla')
>>> log.warning('bla')
>>> log.error('bla')
>>> log.critical('bla')
```

Where these messages end up is configured in `settings_default.py` and you can as usual override the default in the node's own `settings.py`. For example, you set the location and name of the log-file like this:

```
LOGGING['handlers']['logfile']['filename'] = '/path/to/yourlog.log'
```

Note: Critical errors (using `log.critical()`) are sent to the configured admin email address. You need to supply a valid address and make sure your server can send emails.

In the future, messages issued with `log.debug()` will not be logged, when the global flag `DEBUG=False`.

For more information see <https://docs.djangoproject.com/en/1.3/topics/logging/>

OBTAIN AND USE A VIRTUAL MACHINE WITH THE NODESOFTWARE

9.1 About VirtualBox

VirtualBox is a software that allows you to run an operating system inside a “virtual machine” (VM). You need to download the software from <http://www.virtualbox.org/> and install it on your computer which becomes the *host* for the VM.

9.2 The virtual harddisk

The file can be downloaded from <http://vamdc.tmy.se/files/VAMDCnode.vdi.bz2> (550MB) Last update: March 10, 2011.

Unpack it (*bunzip*) and save it wherever it pleases you (unpacked size is over 2GB). The default location on Linux hosts is `~/VirtualBox/HardDisks/`.

9.3 Setting up the VM

After installing VirtualBox you start it and click “New” for setting up a new VM. Choose Linux/Debian as operating system and a memory size that comfortably fits into your RAM. When asked to create a virtual hard disk, chose the one you downloaded instead.

When you finished the setup, you can click “Start” to boot the VM. At the end of the boot process, you will see a login prompt. Use `vamdc` as username and `V@mdc` as password.

9.4 Once inside the VM

9.4.1 Passwords

The first thing to do is change the password by typing `passwd`. The user can execute commands with root-privileges by prepending `sudo`. Use `sudo passwd` to change the root password (which is also `V@mdc` from the start).

9.4.2 Network

Check that the VM has access to the network by trying to *ssh* to another machine or by running `sudo ifconfig` and checking that interface *eth0* has an IP-address assigned to it.

Note: After copying a VM, the operating system still remembers the old host's network card and will likely give a new name (*eth1*) to the current network interface. You can fix this by removing the content of the file `/etc/udev/rules.d/70-persistent-net.rules` and reboot (type `sudo reboot`); alternatively by editing `/etc/network/interfaces` and running `/etc/init.d/networking restart`.

9.4.3 Install system updates

Keep the system up to date with:

```
$ sudo apt-get update
$ sudo apt-get upgrade
```

9.4.4 The Node Software

You find the Node Software in the home directory of the user *vamdc*:

```
$ cd
$ cd NodeSoftware
$ # alternatively: cd $VAMDCROOT
$ # this environment variable is pre-set.
```

This is a version control repository and you can update the software by:

```
$ git pull upstream
```

For more information on how to use this, please read *Collaborating with git and GitHub* in the previous section. (You can connect the existing repository to your GitHub account with `git add remote origin <YourGitHubRepoURL>` after forking the main repository there.)

Now you should be all set to continue with the *Step by step guide to a new VAMDC node*.

9.4.5 Deployed node

In the VM, both *nginx* and *gunicorn* are installed, as described in *Deployment of your node*. There is a symbolic link `NodeSoftware/nodes/RunningNode` which points to the `ExampleNode`. Once you made your copy of the `ExampleNode` (see *Step by step guide to a new VAMDC node*), point `RunningNode` to your own instead since *nginx* uses `NodeSoftware/nodes/RunningNode/nginx.conf` for its config. Don't forget to restart *nginx* with `service nginx restart`.

In your node directory, you can now run `gunicorn_django -c gunicorn.conf` to start the node workers and should have a running node. To access it from outside the VM, you must probably tweak the network setup between the VM and your host computer.

9.4.6 MySQL

MySQL server 5.1 is installed in the VM. You get a MySQL-prompt with:

```
$ mysql -u root -p
```

The password is once more V@mdc. From this prompt you can create new databases and set the access rights to match the ones from your node's `settings.py`.

Typical commands would be:

```
mysql> CREATE DATABASE yourDBname CHARACTER SET utf8;
mysql> GRANT ALL PRIVILEGES ON yourDBname.* TO YourUser@localhost IDENTIFIED BY "reeH5ohm";
mysql> flush privileges;
```

ADDITIONAL TOPICS

10.1 Setting the related name of a field

When you have a *ForeignKey* called *key1* in a *ModelB* which points *ModelA*, the fields from *ModelA* become accessible by *b.key1.fieldFromModelA* in a selection *b* of *ModelB*. This is using the *ForeignKey* in **forward direction**.

Django also automatically adds a field to *ModelA* that contains all the instances of *ModelB* that point to a specific instance *a* of *ModelA*. This field is by default called as the referenced model plus *_set*. So *a.modelb_set* would hold all the *ModelBs* that reference *a*. This is using the *ForeignKey* in **inverse direction**.

You can change the name of the inverse field by giving the argument *related_name='bla'* to the definition of the *ForeignKey* in the model. When you have more than one *ForeignKey* from one model to the same other model, you **must** set the *related_name* because the automatic naming cannot give the same name twice.

A typical example for this are the upper and lower states for a transition where it makes sense to have two *ForeignKeys* in the *Transition* model, e.g. called *upstate* and *lostate*, each pointing to an entry in the *State* model. Now one sets the *related_names* of these *ForeignKeys* to something like *'transitions_with_this_upstate'* and *'transitions_with_this_lostate'* respectively. Thereby, for any state *s* the transitions that have *s* as upper state can be retrieved by *s.transitions_with_this_upstate*.

10.2 Inserting custom XML into the generator

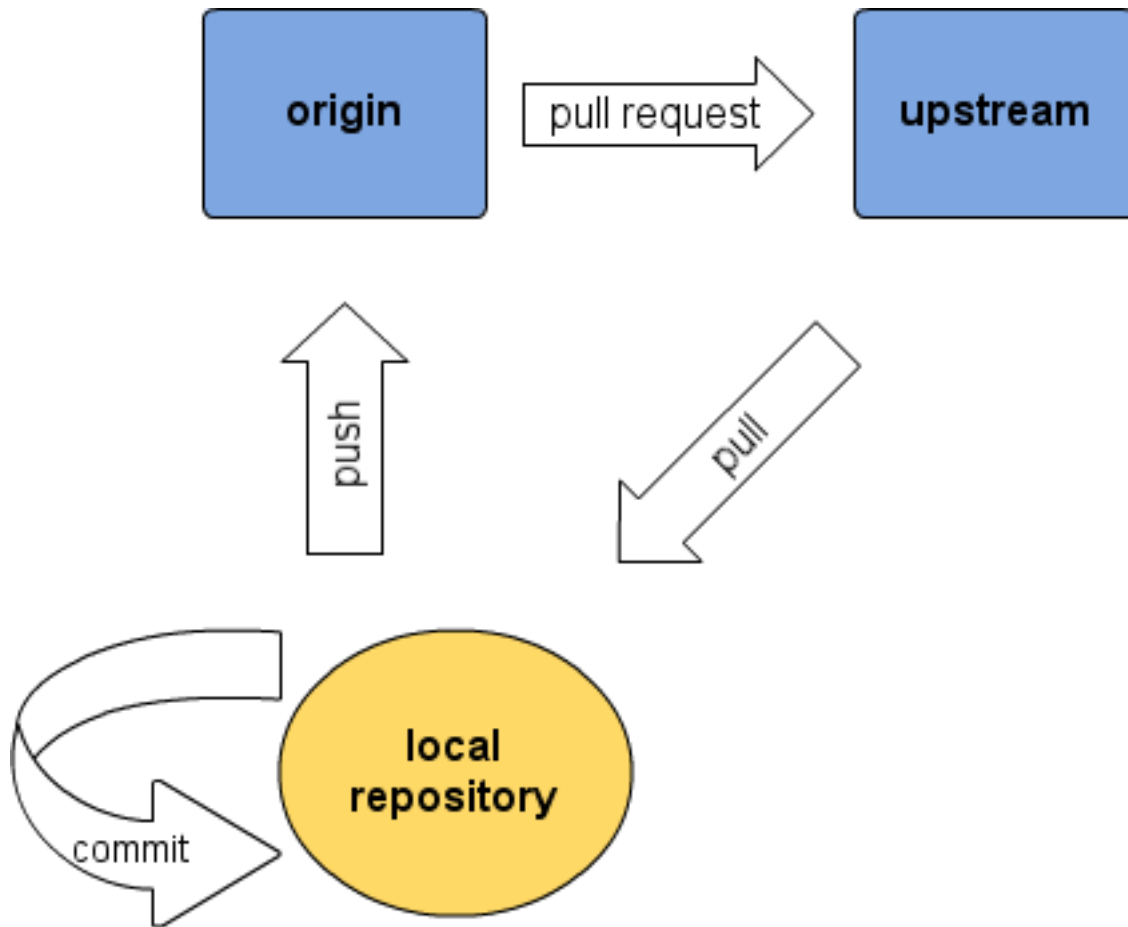
There can arise situations where it might be easier for a node to create a piece of XML itself than filling the *Returnable* and letting the generator handle this. This is allowed and the generator checks every time it loops over an object, if the loop variable, e.g. *AtomState* has an attribute called *XML*. If so, it returns *AtomState.XML()* instead of trying to extract the values from the *Returnable* for the current block of *XSAMS*. Note the *execution* of *.XML()* which means that this needs to be coded as a function/method in your model, not as an attribute.

10.3 Collaborating with git and GitHub

Git is a decentralized version control system (<http://git-scm.com/>). This means among other things that:

- Each checked out copy of the code has the full version history.
- There is no central repository, all repositories (“repos”) are equal (but some *can* be made more equal than others, as we’ll see below).
- Commits happen locally into your working repo, no network connection needed.
- Repos are updated and synced with each other by pushing and pulling commits back and forth between them.
- There are web-platforms that offer free web-repositories which facilitates syncing and merging. We’ll use *GitHub* (<http://www.github.com/>).

The setup that we want looks like this:



- The **local repository** (also known as your “working copy”) is your own workspace. This is where you do all your work. It offers you full local version control without necessarily having to upload the changes anywhere. We’ll get to how you create your local repo in a minute.
- Your **origin** is an online version of your repository, stored online at GitHub. When you want to sync the two you need to *push* your latest local changes to origin. Once online, others will also be able to see the changes.
- **Upstream** is a unique repository that serves as an online code “central” managed by VAMDC. It too is hosted on GitHub. Upstream serves as a convenient way to update your distribution; you should regularly *pull* the latest changes into your local repo to stay updated. Conversely, if you want your own changes to be incorporated into the central distribution you can send a *pull request* to upstream. The relevant commit(s) in your **origin** repo will be reviewed and will, if accepted, be merged into upstream so that others will get the changes next time they do a pull.
- You can certainly have **several local repositories**, e.g. one on your laptop, one on your desktop and one on the server where the node runs. You then use the online **origin** repository to keep them in sync. For example: You work from your laptop and commit your changes locally. You then push them to your origin repository. Next all you need to do is to tell your other local repos to pull from origin and they will all be synced.

Now enough with theory, let’s do this in practice. To create your own repositories (origin and local) do the following:

- Go to <http://github.com> and make an account. This includes that you (create and) upload an ssh-key to be able to pull and push securely and without typing your password all the time. Simply follow the instructions on GitHub.
- Visit the repository at <https://github.com/VAMDC/NodeSoftware> and click “fork” in the upper right corner. This will make a copy of the original repository under your account. This is your **origin** (see above). For

more information on forking, you can read <http://help.github.com/forking/>.

- Github will give you instructions on how to *clone* your origin to your own computer, thereby creating a local repo, your **local repository**, aka your “working copy”.
- You can repeat the cloning on as many machines as you see fit.
- Tell your local repos where **upstream** is by running the following command in each of them: `git remote add upstream git://github.com/VAMDC/NodeSoftware.git`

Now that you are all set, a typical working session may look like this:

```
$ cd $VAMDCROOT # got to your local repo
$ git status # should tell you you have a clean tree and are on the branch "master"
$ git pull origin # pull from your origin, in case you pushed things there from another
$ git pull upstream # fetch the latest from upstream and merge it with your tree.
$ git log # read the commit log about what is new.
$ ... # edit your files
$ git status # review which files have changed
$ git diff # review details of your changes
$ git diff <filename> # see changes in one file only
$ git add <filename> # add a file to be committed with the next commit, e.g. a new file
$ git commit -a -m "message" # commit all changed files. ALWAYS check the status before you use
$ git commit -m "message" <filenames> # commit, but include only the named files in the commit
$ ... # more edits, more commits. until, at the end of day:
$ git status # also tells you how many commits you are ahead of your origin
$ git push # push all commits to your origin, also the new ones that came from
```

Note: There are several graphical user interfaces available for git that will facilitate overview and some operations for the less command-line adept. Commonly used ones for Linux are *gitk* and *gitg*. Good editors also integrate with git so that you can handle the version control from within the editor.

After you pushed your work to your origin, you can go to the *GitHub* website and send a *pull request* to the upstream repository, if you want your changes to be propagated to everybody else. We will then look at your commits and merge them.

A few dos and don'ts that are worthwhile to keep in mind with git:

- Do commit often. It goes instantly.
- Pull and push less often, but often enough. You certainly want to pull from upstream before making changes, since you otherwise might work on outdated versions of files which will result in conflicts later. You also do not want to sit on your local commits for too long but push them frequently instead.
- Never pull into a dirty tree (i.e. one that has uncommitted changes). Commit first, then pull. Alternatively read *git help stash*.
- Do *not* commit data files that you have put in your node directory. (check `git status` on what will be committed before you use `git commit -a`.)
- *Git* trusts you know what you are doing. It will allow you to do stupid things, too.
- Don't panic. Yes, *git* may have a comparably steep learning curve, but it is a powerful tool and all problems can be resolved.

10.3.1 Situations that commonly arise and how to solve them

Merge conflicts. When you pull from Upstream into your repo, other's changes are merged with yours. It might however happen that someone else has changed the same line in the same file as you have in onw of your own commits, which results in a merge conflict. The pull commands warns you about this and *git status* shows the file in question as “both modified”. The file itself contains both versions of the conflicting lines, clearly marked. Edit the file so that only one version remains and remove the markings. Then you simply commit the file (and push).

Undo a commit. To undo a commit means exactly that, *not* that any of the files change. For example, undoing the last commit leaves you with as much uncommitted changes as you had before your last commit. None of your edits is reversed. Undoing commits is practical e.g. when you have committed too many things at once or unwanted files; or when you want to split one commit into several. You undo a commit with `git reset --soft <REF>` where `<REF>` is the commit that should be resetted to (i.e. the next-to-last one, if you want to undo your last commit). Common values for `<REF>` include:

- `HEAD^` - this is the next-to-last
- `HEAD^^` - the one before the next-to-last.
- `HEAD~5` - five commits back
- `111521cb9d3771e636f5f053d3d1048aa7c8852f` - each commit has a long hash number that uniquely identifies it. They can be seen in `git log` and you can give the hash number of the commit that you want to reset to to `git reset`.

Revert to an earlier version. If you want to *throw away* your edits since a certain commit, you use `git reset --hard`. For example, to revert all files to the state that they were in at the last commit (throw away uncommitted changes), you do `git reset --hard HEAD`. Similarly to the soft reset, you can also specify earlier commits that you want to reset to.

Look at an earlier version. You can check out any earlier version of any file at any time. For example, `git checkout "master@{1 month ago}" <filename>` will give you the version of the file `<filename>` from a month ago. To go back to the latest, you do `git checkout master <filename>` ("master" is the name of the default branch where all your commits are). Note that the last command can also be used to throw away uncommitted changes in a specific file - a more gentle way than the reset described above.

You can also skip the `<filename>` to check out an earlier version of the whole repo (`git checkout master` brings you back to the latest). Instead of "master@{1 month ago}" you can use any of the `<REF>` mentioned above, or have a look at http://book.git-scm.com/4_git_treeishes.html.

Make a branch. Read `git help branch` for this.

10.3.2 Commit guidelines

One thing at a time. Please commit often and only include things in one commit that logically belong together. For example, changes to your node and changes to the common library should not be in the same commit but committed separately.

Meaningful commit messages. This goes together with the previous: If you cannot meaningfully summarize the changes you want to commit in one or two lines, your commit is likely to be too large. Try to make the log messages meaningful!

Good code. Please try to avoid spaghetti-code, write modular, and follow <http://www.python.org/dev/peps/pep-0008/>

Pull first. Before you send a pull request, please make sure that you have pulled from upstream. This will make the merging of your code easier, since it will be you who needs to resolve potential conflicts before you push to your origin again.

The admin of *upstream* (aka the writer of these lines) might be bribed and/or convinced to turn a blind eye on violations against any of the above points, but he will be very happy if you try to follow them.

10.4 Adding more views or apps to your node

tbw

10.5 The Django admin interface

tbw

10.6 Handling advanced queries

tbw

10.7 Using a custom model method for filling a Returnable

tbw

KNOWN LIMITATIONS

In general, the NodeSoftware tries to be forgiving with faulty input data from the nodes' databases and will do its best to return a valid and complete XML document. However it relies on the content of the connected database and the connection to the schema via the models and dictionaries. Errors in these cannot be compensated by the software itself and can result in invalid output data. All nodes are encouraged to check the validity of their XML output against the current XSAMS, for example with the help of the TAPValidator application.

The NodeSoftware does and will not offer the full possibilities of the XSAMS since choices and simplifications have to be made in the implementation. These deliberate limitations include:

- Treating isotopes and ions of atoms as different species, repeating the element information instead of nesting several ions within each isotope, and nesting the ions within each element.
- Only allowing one set of quantum numbers per atomic or molecular state. If a node wishes to return several different descriptions of the quantum numbers per state, this needs to be implemented in a custom fashion for this node.
- Only one set of line broadening parameters per transition and per type (instrumental, natural, pressure, doppler) is allowed at this time. The next release of the software will include the possibility to give several pressure-broadenings per transition.

Tools for handling more advanced queries and for treating certain Restrictables as special cases in a node's query-function are lacking and will be improved in the next version.

The SQL-parser does currently not support advanced nested queries with several levels of brackets. Also the IN-operator is as of now unsupported. This will be amended in the next release.

A full list of outstanding issues is available at the development repository at <https://github.com/VAMDC/NodeSoftware/issues> where anybody is welcome to file bugs or wishlist-items.

BUGS AND CONTACT

12.1 Report a bug

The NodeSoftware is in active development and there are some rough edges still. We very much appreciate your feedback and complaints are usually resolved quickly.

Please file an issue at <https://github.com/VAMDC/NodeSoftware/issues>

This can be used for bugs on both, the software itself and the documentation.

12.2 Contact information

You can write to the VAMDC developers email list: `vamdc.developer <AT> sympa.obspm.fr`

THE CODE

You can download the NodeSoftware at these locations:

- Release 11.5: <http://www.vamdc.eu/downloads/NodeSoftware-v11.5r1.tar.gz>
- Release branch: <https://github.com/VAMDC/NodeSoftware/tarball/release>
- Latest development branch: <https://github.com/VAMDC/NodeSoftware/tarball/master>

Please see *Upgrading*.

The development repository resides at <https://github.com/VAMDC/NodeSoftware> and you are welcome to use the version control software *git* to check out your own copy. This takes a few more minutes to set up but has the benefit of facilitating collaboration. After all, you might makes changes or extend the code for your needs and we would like to include your improvements into the main repository. Read more about this at *Collaborating with git and GitHub*.

13.1 Source code documentation

The following is the automatically generated documentation from the source code. It lists and describes all functions, classes etc.

13.2 The VAMDC-TAP service library

13.2.1 Tapservice Documentation

This page contains the Tapservice Package documentation.

The `sqlparse` Module

```
vamdctap.sqlparse.setupSQLparser ()  
vamdctap.sqlparse.singleWhere (w, RESTRICTABLES)  
vamdctap.sqlparse.where2q (ws, RESTRICTABLES)
```

The `generators` Module

```
vamdctap.generators.GetValue (name, **kwargs)  
    the function that gets a value out of the query set, using the global name and the node-specific dictionary.
```

`vamdctap.generators.Xsams` (*HeaderInfo=None, Sources=None, Methods=None, Functions=None, Environments=None, Atoms=None, Molecules=None, Solids=None, Particles=None, CollTrans=None, RadTrans=None, RadCross=None, NonRadTrans=None*)

The main generator function of XSAMS. This one calls all the sub-generators above. It takes the query sets that the node's `setupResult()` has constructed as arguments with given names. This function is to be passed to the HTTP-respose object directly and not to be looped over beforehand.

`vamdctap.generators.XsamsAtoms` (*Atoms*)

Generator (yield) for the main block of XSAMS for the atoms, with an inner loop for the states. The `QuerySet` that comes in needs to have a nested `QuerySet` called `States` attached to each entry in `Atoms`.

`vamdctap.generators.XsamsCollTrans` (*CollTrans*)

Collisional transitions. `QuerySets` and nested `querysets`: # `CollTran` # `CollTran.Reactants` # `CollTran.IntermediateStates` # `CollTran.Products` # `CollTran.DataSets` # `DataSet.FitData` # `FitData.Arguments` # `FitData.Parameters` # `DataSet.TabulatedData`

Matching loop variables to use: # `CollTran` # `CollTranReactant` # `CollTranIntermediateState` # `CollTranProduct` # `CollTranDataSet` # `CollTranFitData` # `CollTranFitDataArgument` # `CollTranFitDataParameter` # `CollTranTabulatedData`

`vamdctap.generators.XsamsEnvironments` (*Environments*)

`vamdctap.generators.XsamsFunctions` (*Functions*)

Generator for the `Functions` tag

`vamdctap.generators.XsamsMCSBuild` (*Molecule*)

Generator for the `MolecularChemicalSpecies`

`vamdctap.generators.XsamsMSBuild` (*MoleculeState*)

Generator for `MolecularState` tag

`vamdctap.generators.XsamsMethods` (*Methods*)

Generator for the `methods` block of XSAMS

`vamdctap.generators.XsamsMolecules` (*Molecules*)

Generator for `Molecules` tag

`vamdctap.generators.XsamsNonRadTrans` (*NonRadTrans*)

non-radiative transitions

`vamdctap.generators.XsamsParticles` (*Particles*)

Generator for `Particles` tag.

`vamdctap.generators.XsamsRadCross` (*RadCross*)

for the `Radiative/CrossSection` part

`querysets` and nested `querysets`:

RadCros

RadCros.BandAssignments

BandAssignment.Modes `Mode.DeltaVs`

loop variables:

RadCros

RadCrosBandAssignment

RadCrosBandAssignmentMode `RadCrosBandAssignmentModeDeltaV`

`vamdctap.generators.XsamsRadTranBroadening` (*G*)

helper function for line broadening, called from `RadTrans`

allowed names are: `pressure`, `instrument`, `doppler`, `natural`

`vamdctap.generators.XsamsRadTranShifting` (*RadTran, G*)

Shifting type

`vamdctap.generators.XsamsRadTrans` (*RadTrans*)
Generator for the XSAMS radiative transitions.

`vamdctap.generators.XsamsSolids` (*Solids*)
Generator for Solids tag

`vamdctap.generators.XsamsSources` (*Sources*)
Create the Source tag structure (a bibtex entry)

`vamdctap.generators.checkXML` (*obj*)
If the queryset has an XML method, use that and skip the hard-coded implementation.

`vamdctap.generators.countReturnables` (*regex*)
count how often a certain matches the keys of the returnables

`vamdctap.generators.embedhtml` (*transitions, totalcount=None*)
Embed html

`vamdctap.generators.escape` (*s*)

`vamdctap.generators.generatorError` (*where*)

`vamdctap.generators.isiterable` (*obj*)

`vamdctap.generators.makeArgumentType` (*tagname, keyword, G*)
Build ArgumentType

`vamdctap.generators.makeAtomComponent` (*Atom, G*)
This constructs the Atomic Component structure.

Atom - the current Atom queryset G - the shortcut to the GetValue function

`vamdctap.generators.makeBroadeningType` (*G, name='Natural'*)
Create the Broadening tag

`vamdctap.generators.makeCaseQNs` (*G*)
return the Case and the QNs

`vamdctap.generators.makeDataFuncType` (*tagname, keyword, Parameter, G*)
Build the DataFuncType.

`vamdctap.generators.makeDataSeriesType` (*tagname, keyword, G*)
Creates the dataseries type

`vamdctap.generators.makeDataType` (*tagname, keyword, G, extraAttr=None, extraElem=None*)
This is for treating the case where a keyword corresponds to a DataType in the schema which can have units, comment, sources etc. The dictionary-suffixes are appended and the values retrieved. If the sources is iterable, it is looped over.

#This extends the PrimaryType with some often-seen arguments.

`vamdctap.generators.makeNamedDataType` (*tagname, keyword, G*)
Similar to makeDataType above, but allows the result of G() to be iterable and adds the name-attribute. If the corresponding refs etc are not iterable, they are replicated for each tag.

`vamdctap.generators.makePartitionfunc` (*keyword, G*)
Create the Partionfunction tag element.

`vamdctap.generators.makePrimaryType` (*tagname, keyword, G, extraAttr=None*)
Build the Primary-type base tags. Note that this method does NOT close the tag, </tagname> must be added manually by the calling function.

extraAttr is a dictionary of attributes-value pairs to add to the tag.

`vamdctap.generators.makeShellType` (*tag, keyword, G*)
Creates the Atom shell type.

`vamdctap.generators.makeSourceRefs` (*refs*)
Create a SourceRef tag entry

`vamdctap.generators.makeTermType` (*tag, keyword, G*)

Construct the Term xsams structure.

This version is more generic than XsamsTerm function and don't enforce LS/JK/LK to be exclusive to one another (as dictated by current version of xsams schema)

`vamdctap.generators.makeiter` (*obj*)

Return an iterable, no matter what

`vamdctap.generators.makeloop` (*keyword, G, *args*)

Creates a nested list of lists. All arguments should be valid dictionary keywords and will be fed to G. They are expected to return iterables of equal lengths. The generator yields a list of current element of each argument-list in order, so one can do e.g.

for name, unit in makeloop('TabulatedData', G, 'Name', 'Unit'): ...

`vamdctap.generators.parityLabel` (*parity*)

XSAMS whats this as strings "odd" or "even", not numerical

`vamdctap.generators.sources2votable` (*sources*)

Sources to VO

`vamdctap.generators.states2votable` (*states*)

States to VO

`vamdctap.generators.transitions2embedhtml` (*transs, count*)

Converting Transition to html

`vamdctap.generators.transitions2votable` (*transs, count*)

Transition to VO

`vamdctap.generators.votable` (*transitions, states, sources, totalcount=None*)

VO base definition

The views Module

class `vamdctap.views.TAPQUERY` (*data*)

Bases: object

This class holds the query, does some validation and triggers the SQL parser.

validate ()

`vamdctap.views.addHeaders` (*headers, response*)

`vamdctap.views.async` (*request*)

`vamdctap.views.availability` (*request*)

`vamdctap.views.availabilityXsl` (*request*)

`vamdctap.views.capabilities` (*request*)

`vamdctap.views.capabilitiesXsd` (*request*)

`vamdctap.views.capabilitiesXsl` (*request*)

`vamdctap.views.cleandict` (*dict*)

throw out keys where the value is ''

`vamdctap.views.getBaseURL` (*request*)

`vamdctap.views.index` (*request*)

`vamdctap.views.randStr` (*n*)

`vamdctap.views.sync` (*request*)

`vamdctap.views.tables` (*request*)

`vamdctap.views.tablesXsd` (*request*)

`vamdctap.views.tapNotFoundError` (*request*)

`vamdctap.views.tapServerError` (*request=None, status=500, errmsg=''*)

13.3 The import tool

13.3.1 Imptools Documentation

This page contains the Imptools Module documentation.

The `imptools` Module

This program implements a database importer that reads from ascii-input files to the django database. It's generic and is controlled from a mapping file.

class `imptools.rewrite.MappingFile` (*filepath, headblocks, commentchar, blockoffset, blockstep, errblock, startblock=None, endblock='n'*)

Bases: `object`

This class implements an object that represents an open file from which one can read blocks. The object keeps track of its own block-step speed and will return lines as defined by this speed. E.g. for a line-step speed of 0.5, it will return the same line twice in a row whereas for a step speed of 2, will return every second block etc.

If `endblock` is `\n` (default), the block will infact represent a line.

block_generator (*fileobj, startblock=None, endblock='n'*)
generator, stepping through blocks

readblock ()

Return a block from the file.

This method understand both slower block stepping ($0 < \text{blockstep} < 1$) and faster (> 1)

`imptools.rewrite.ftime` (*t0, t1*)
formats time to nice format.

`imptools.rewrite.get_value` (*linedata, column_dict*)
Process one line of data. Linedata is a tuple that always starts with the raw string for the line. The function with its arguments is read from the `column_dict` and applied to the `linedata`. The result is returned, after checking for the NULL value.

`imptools.rewrite.is_iter` (*iterable*)

`imptools.rewrite.log_trace` (*e, info=''*)
Intended to be called from inside a traceback exception with the exception object as first argument. Captures the latest traceback.

`imptools.rewrite.make_outfile` (*file_dict, global_debug=False*)
Process one file definition from a config dictionary by processing the file name stored in it and parse it according to the mapping.

`file_dict` - config dictionary representing one input file structure (for an example see e.g. `mapping_vald3.py`)

`imptools.rewrite.parse_mapping` (*mapping, debug=False*)
Step through a list of mappings describing the relation between (usually ascii-)files and django database fields. This should ideally not have to be changed for different database types.

`imptools.rewrite.read_mapping` (*fname*)
Read the config dictionary from a file. Note: very unsafe, since the content gets executed. Have a look at `creatcfg()` to see how it should look like.

`imptools.rewrite.validate_mapping` (*mapping*)

Check the mapping definition to make sure it contains a structure suitable for the parser.

13.3.2 Linefuncs Documentation

This page contains the Linefuncs Module documentation. See also *How to get your data into the database*.

The linefuncs Module

Line functions are helper functions available to use in the mapping file.

All importable line functions take 'linedata' as a first argument. This can be either a string (the current active line), or a list of strings in the case of lines from many files being read simultaneously. In the latter case, the argument 'filenum' is used to select the correct line. For linefuncs accepting a single line, this selection must be made in the mapping dictionary.

Example of call from mapping dictionary:

```
{'cname' ['whatever_field_name', ] 'cbyte' : (charrange, 56, 58)}
```

`imptools.linefuncs.bySepNr` (*linedata, number, sep=''*, *'*)

`imptools.linefuncs.bySepNr2` (*linedata, number, sep=''*, *'*, *filenum=0*)

Split a text line by sep argument and return the number:ed split section

Inputs: linedata (str or iterable) - current line(s) to operate on number (int) - nth section, separated by sep (str) - a separator to split by filenum (int) - optional selection of line if linedata is an iterable

`imptools.linefuncs.charrange` (*linedata, start, end, filenum=0*)

Cut out part of a line of texts based on indices.

Inputs: linedata (str or iterable) - current line(s) to operate on start, end (int) - beginning and end indices of the line filenum (int) - optional selection of line if linedata is an iterable

`imptools.linefuncs.charrange2int` (*linedata, start, end, filenum=0*)

Cut out part of a line based on indices, return as integer

Inputs: linedata (str or iterable) - current line(s) to operate on start, end (int) - beginning and end indices of the line filenum (int) - optional selection of line if linedata is an iterable

`imptools.linefuncs.constant` (*linedata, value*)

`imptools.linefuncs.get_accur` (*linedata, range1, range2*)

extract accuracy

`imptools.linefuncs.get_alphawaals` (*linedata, sep1, sep2*)

extract alpha - van der waal value

`imptools.linefuncs.get_gammawaals` (*linedata, sep1, sep2*)

extract gamma - van der waal value

`imptools.linefuncs.get_publications` (*linedata*)

extract publication data. This returns a list since it is for a multi-reference.

`imptools.linefuncs.get_sigmawaals` (*linedata, sep1, sep2*)

extract sigma - van der waal value

`imptools.linefuncs.get_srcfile_ref` (*linedata, sep1, sep2*)

extract srcfile reference

`imptools.linefuncs.get_term_val` (*linedata, n*)

extract configurations from term file

`imptools.linefuncs.is_iter` (*iterable*)

Helper function

Checks if the given argument is iterable or not, i.e. if it is a list or tuple. Strings are not considered iterable by this function.

`imptools.linefuncs.lineSplit` (*linedata*, *splitsep*=' ', *filenum*=0)

Splits a line by *splitsep*, returns a list. The main use for this method is creating a many-to-many reference.

Inputs: *linedata* (str or iterable) - current line(s) to operate on
splitsep (str) - string to split by
filenum (int) - optional selection of line if *linedata* is an iterable

Returns a list!

`imptools.linefuncs.merge_cols` (*linedata*, **ranges*)

Merges data from several columns into one, separating them with '-'. *ranges* are any number of tuples (*indexstart*, *indexend*) defining the columns.

`imptools.linefuncs.merge_cols_by_sep` (*linedata*, **sepNr*)

Merges data from several columns (separated by ;) into one, separating them with '-'. *sepNr* are the nth position of the file, separated by 'sep'. Assumes a single line input.

PYTHON MODULE INDEX

i

`imptools.linefuncs`, 49
`imptools.rewrite`, 48

v

`vamdctap.generators`, 44
`vamdctap.sqlparse`, 44
`vamdctap.views`, 47

INDEX

A

addHeaders() (in module vamdctap.views), 47
async() (in module vamdctap.views), 47
availability() (in module vamdctap.views), 47
availabilityXsl() (in module vamdctap.views), 47

B

block_generator() (imptools.rewrite.MappingFile
method), 48
bySepNr() (in module imptools.linefuncs), 49
bySepNr2() (in module imptools.linefuncs), 49

C

capabilities() (in module vamdctap.views), 47
capabilitiesXsd() (in module vamdctap.views), 47
capabilitiesXsl() (in module vamdctap.views), 47
charrange() (in module imptools.linefuncs), 49
charrange2int() (in module imptools.linefuncs), 49
checkXML() (in module vamdctap.generators), 46
cleandict() (in module vamdctap.views), 47
constant() (in module imptools.linefuncs), 49
countReturnables() (in module vamdctap.generators),
46

E

embedhtml() (in module vamdctap.generators), 46
escape() (in module vamdctap.generators), 46

F

ftime() (in module imptools.rewrite), 48

G

generatorError() (in module vamdctap.generators), 46
get_accur() (in module imptools.linefuncs), 49
get_alphawaals() (in module imptools.linefuncs), 49
get_gammawaals() (in module imptools.linefuncs), 49
get_publications() (in module imptools.linefuncs), 49
get_sigmawaals() (in module imptools.linefuncs), 49
get_srcfile_ref() (in module imptools.linefuncs), 49
get_term_val() (in module imptools.linefuncs), 49
get_value() (in module imptools.rewrite), 48
getBaseURL() (in module vamdctap.views), 47
GetValue() (in module vamdctap.generators), 44

I

imptools.linefuncs (module), 49
imptools.rewrite (module), 48
index() (in module vamdctap.views), 47
is_iter() (in module imptools.linefuncs), 49
is_iter() (in module imptools.rewrite), 48
isiterable() (in module vamdctap.generators), 46

L

lineSplit() (in module imptools.linefuncs), 50
log_trace() (in module imptools.rewrite), 48

M

make_outfile() (in module imptools.rewrite), 48
makeArgumentType() (in module vamdc-
tap.generators), 46
makeAtomComponent() (in module vamdc-
tap.generators), 46
makeBroadeningType() (in module vamdc-
tap.generators), 46
makeCaseQNs() (in module vamdctap.generators), 46
makeDataFuncType() (in module vamdc-
tap.generators), 46
makeDataSeriesType() (in module vamdc-
tap.generators), 46
makeDataType() (in module vamdctap.generators), 46
makeeiter() (in module vamdctap.generators), 47
makeloop() (in module vamdctap.generators), 47
makeNamedDataType() (in module vamdc-
tap.generators), 46
makePartitionfunc() (in module vamdctap.generators),
46
makePrimaryType() (in module vamdctap.generators),
46
makeShellType() (in module vamdctap.generators), 46
makeSourceRefs() (in module vamdctap.generators),
46
makeTermType() (in module vamdctap.generators), 46
MappingFile (class in imptools.rewrite), 48
merge_cols() (in module imptools.linefuncs), 50
merge_cols_by_sep() (in module imptools.linefuncs),
50

P

parityLabel() (in module vamdctap.generators), 47

parse_mapping() (in module imptools.rewrite), 48

R

randStr() (in module vamdctap.views), 47

read_mapping() (in module imptools.rewrite), 48

readblock() (imptools.rewrite.MappingFile method), 48

S

setupSQLparser() (in module vamdctap.sqlparse), 44

singleWhere() (in module vamdctap.sqlparse), 44

sources2votable() (in module vamdctap.generators), 47

states2votable() (in module vamdctap.generators), 47

sync() (in module vamdctap.views), 47

T

tables() (in module vamdctap.views), 47

tablesXsd() (in module vamdctap.views), 47

tapNotFoundError() (in module vamdctap.views), 48

TAPQUERY (class in vamdctap.views), 47

tapServerError() (in module vamdctap.views), 48

transitions2embedhtml() (in module vamdc-
tap.generators), 47

transitions2votable() (in module vamdctap.generators),
47

V

validate() (vamdctap.views.TAPQUERY method), 47

validate_mapping() (in module imptools.rewrite), 48

vamdctap.generators (module), 44

vamdctap.sqlparse (module), 44

vamdctap.views (module), 47

votable() (in module vamdctap.generators), 47

W

where2q() (in module vamdctap.sqlparse), 44

X

Xsams() (in module vamdctap.generators), 44

XsamsAtoms() (in module vamdctap.generators), 45

XsamsCollTrans() (in module vamdctap.generators), 45

XsamsEnvironments() (in module vamdc-
tap.generators), 45

XsamsFunctions() (in module vamdctap.generators), 45

XsamsMCSBuild() (in module vamdctap.generators),
45

XsamsMethods() (in module vamdctap.generators), 45

XsamsMolecules() (in module vamdctap.generators),
45

XsamsMSBuild() (in module vamdctap.generators), 45

XsamsNonRadTrans() (in module vamdc-
tap.generators), 45

XsamsParticles() (in module vamdctap.generators), 45

XsamsRadCross() (in module vamdctap.generators), 45

XsamsRadTranBroadening() (in module vamdc-
tap.generators), 45

XsamsRadTrans() (in module vamdctap.generators), 45

XsamsRadTranShifting() (in module vamdc-
tap.generators), 45

XsamsSolids() (in module vamdctap.generators), 46

XsamsSources() (in module vamdctap.generators), 46