

Writing bibliographic tools with
pybliographer

Frédéric Gobry

February 15, 2006

Contents

1	Introduction	2
1.1	Basic concepts	2
1.1.1	The database schema	3
1.1.2	Taxonomies	3
1.1.3	Result sets	4
1.1.4	Views	4
1.2	Manipulating data	4
1.2.1	Loading and saving	4
1.2.2	Using the registry	5
1.2.3	Updating records	5
1.2.4	Sorting	6
1.2.5	Searching	6
1.3	Importing and exporting	7
1.4	Citation formatting	7
2	Extending <i>pybliographer</i>	9
2.1	Specializing a parser	9

Chapter 1

Introduction

pybliographer is a developer-oriented framework for manipulating bibliographic data. It is written in *python*¹, and uses extensively the dynamic nature of this language.

pybliographer does not try to define another standard format for bibliographic data, nor does it solely rely on a single existing standards. Standards are important in order to allow for interoperability and durability. Unfortunately, real-world data often contain a great number of mistakes, or reflect certain local conventions. *pybliographer* is on the *pragmatic* side of considering these issues as part of its business: most of the parsing tasks can be easily overridden and specialized in order to *fit the code to the data*, and not the other way around.

1.1 Basic concepts

pybliographer deals with sets of `Records`, stored in a so-called `Database`. This database can be actually implemented on top of different systems. Two are available today, one based on a single XML file, using a custom XML dialect, the other based on Berkeley DB², a very efficient database system.

Each record represents an elementary object you want to describe, and has a number of *attributes*. For instance, if you are describing a book, one attribute will be its *title*, another its *ISBN*, etc. Each of these attributes can contain one or more values, all of the same *type*. To continue the description of our book, we probably have the *author* attribute, which contains as many `Person` values as there are authors for the book. All the values of a given attribute are of the same type.

In some cases, simply having this flat key/value model to describe an object is not enough. *pybliographer* allows, for every value of every attribute, to provide a set of *qualifiers*. These qualifiers are also attributes which can hold one or more values. If my book, or information about the book, is available via the internet, I can provide a *link* attribute, but for each of the actual URLs provided, I might wish to add a *description* qualifier, which will indicate, say, if the URL points to the editor's website, or to a review, etc.

This nesting of objects is best described in figure 1.1.

¹see <http://python.org/>

²see <http://www.sleepycat.com/>

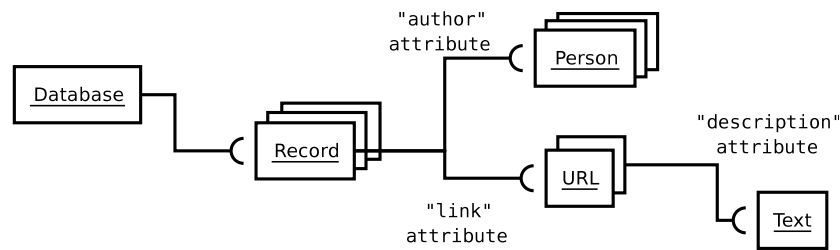


Figure 1.1: Objects manipulated in *pybliographer*

pybliographer comes with a set of defined attribute types, like `Person`, `Text`, `Date`, `ID` (see the `Pyblbio.Attribute` module for a complete list), and can be extended to support your own types.

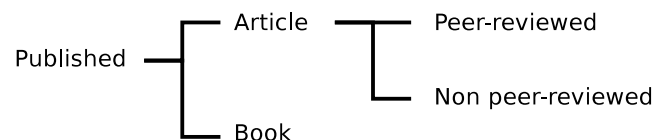
1.1.1 The database schema

Even though attributes are typed, the data model described above is quite flexible. In order for *pybliographer* to help you checking that your records are properly typed, it needs to know the database schema you are using. This schema, usually stored in an XML file with the extension `.sip`, simply lists the known attributes with their type and the qualifiers it allows for its values. Some `.sip` files are distributed with *pybliographer*, and can be seen in the `Pyblbio.RIP` directory.

In addition to validation information, the schema contains human-readable description of the different fields, possibly in several languages, so that it can be automatically extracted by user interfaces to provide up-to-date information.

1.1.2 Taxonomies

Taxonomies can be used as *enumerated values*, say for listing the possible types of a document, or the language in which a text is written. They have however the extra capability of being hierarchical: you can define subcategories of a main category. For instance, imagine a `doctype` taxonomy with the following values:



You can tag an article as `Peer-reviewed`, but you are not required to use the *leaf* values in this tree. In the case you don't know if a publication is reviewed or not, you can use the `Article` tag. Similarly, if you search for all the `Published` documents, you will retrieve all those that have the `Published` tag, but also those that are articles (either peer-reviewed or not), books,...

pybliographer uses the `Pyblbio.Attribute.Txo` object to *represent* a logical value in a given taxonomy. A record can be tagged with this `Txo` object by adding a `Pyblbio.Attribute.TxoItem` value in the corresponding attribute.

Taxonomies can be declared and pre-filled in a database schema, so that any database created from the schema will at least contain the specified taxonomies.

To see how these taxonomies can be further created and modified, please have a look at the `txo` member of a `Database` object, which is an instance of the `Pyblbio.Store.TxoGroup` class.

1.1.3 Result sets

Result sets are used to manipulate an explicit list of records, among all the records kept in a database. They are returned from queries on the database, and can be manipulated by the user. Result sets are somewhat like mathematical sets, as you cannot put duplicate values in them, and they have no default ordering of their elements. You can create result sets via the `rs` attribute of your database, which is an instance of the `Pyblbio.Store.ResultSetStore`.

A special result set is available as `Pyblbio.entries`, and contains at every time **all** the records of the database.

1.1.4 Views

We have seen that result sets are **not** ordered. However, in many cases, one needs to provide the records in a specific order. To do so, you can create a *view* on top of a result set. This view is created by calling the `view` method of the result set, with an `order` parameter being the description of the sort order you wish to have. The module `Pyblbio.Sort` provides elementary constructs to build such a description.

Once the view is created, modifying the corresponding result set leads to updating the view accordingly.

1.2 Manipulating data

This section describes some simple operations you can perform on some subset of a *pybliographer* database.

1.2.1 Loading and saving

The first thing you need to do is of course *actually having* a database available. The following code does the job:

```
from Pyblbio import Store, Schema

schema = Schema.Schema('myschema.sip')
store = Store.get('file')

db = store.dbcreate('mydb.bip', schema)
```

This example relies on the fact that you already have a schema at hand. There are schemas available in the `Pyblio.RIP` directory. It starts by reading the schema. The next step is to select the actual physical store which will hold your database. We choose to store it in a simple XML file, whose canonical extension is `.bip`. The last operation actually creates the database with the specified schema.

Independently of the selected store, it is always possible to *export* a database in the `.sip` format, by calling the `db.xmlwrite(...)` method of the database. Such a file can then be reused later on by using `store.dbimport(...)` instead of `store.dbcreate(...)`.

When you have finished modifying your database, you can call `db.save()` method to ensure that it is properly saved.

Caution: the `bsddb` store for instance is updated at every actual modification, not only when you call the `save` method. Don't rely on it to provide some kind of *rollback* feature.

1.2.2 Using the registry

pybliographer has a mechanism to register known schemas, and specify which import and export filters can properly work with each schema. This mechanism can be used to create our database by asking for a specific schema, as shown below:

```
from Pyblio import Store, Registry

Registry.parse_default()

schema = Registry.getSchema("org.pybliographer/bibtex/0.1")
store = Store.get('file')

db = store.dbcreate('mydb.bip', schema)
```

The registry must be first initialized. Then you can ask for a specific schema, in that case a schema that supports BibTeX databases.

1.2.3 Updating records

The next example will loop over all the records in a database, and add a new author to the list of authors.

```
from Pyblio import Attribute

for record in db.entries.itervalues():
    person = Attribute.Person(last=u"Gobry",
                              first=u"Frédéric")

    record.add('author', person)

    db[record.key] = record

db.save()
```

We use the `itervalues()` iterator to loop over all the records stored in the database. Then, we simply insert a new value in the `author` attribute. The `record.add(...)` method takes care of creating the attribute if it does not exist yet.

One thing not to forget is to store the record back in the database once the modification is performed. Without this step, you might experience weird behavior where some modifications are not properly kept.

We finish by saving the database.

1.2.4 Sorting

To sort records, you create *views* (see section 1.1.4 on page 4). You can of course create multiple views on top of a single result set. In order to sort the whole database, simply create the view on `database.entries` instead of a result set. If you want to sort your database by decreasing year and then by author, you can use a view like that:

```
from Pyblbio.Sort import OrderBy

view = db.entries.view(OrderBy('year', asc=False) &
                       OrderBy('author'))

for record in view.itervalues():
    # do something with the record
    # ...
```

So, sorting constraints can be arbitrarily chained with the `&` operator, and each constraint can be either *ascending* (the default), or *descending*. This defines a very simple *Domain Specific Language*, or DSL for short. Such languages also appear in other part of *pybliographer* (searching, citation formatting), as they are a convenient way to describe complex abstraction without having to reinvent a complete environment.

1.2.5 Searching

To search, you call the `database.query(...)` method. The method takes a query specification as argument, which is constructed with the help of another DSL, similar to the one used for sorting. You have access to a certain number of primitive queries, which are then linked together with the usual boolean operators, as in the following example:

```
from Pyblbio import Query

article = db.txo['doctype'].byname('article')

result = db.query(~ Query.Txo('doctype', article) &
                  Query.AnyWord('laziness'))
```

We first get the taxonomy item corresponding to articles, and we then compose the following query: get all the documents that are *not* articles, and which contain the word *laziness* in any attribute.

1.3 Importing and exporting

As *pybliographer* is not bound to a single data schema, importing and exporting from specific formats (like MARC, BibTeX, Dublin Core,...) cannot be achieved once for all. In order to avoid the need to recreate a BibTeX parser for every database schema invented, *pybliographer* makes a clear separation between *syntactic parsers*, located in `Pybl.io.Parsers.Syntactic` and *semantic parsers*, in `Pybl.io.Parsers.Semantic`. A syntactic parser is only in charge of analyzing or generating a file format, without any assumption regarding the meaning of the fields it reads. These syntactic parsers are then reused by the semantic code, which relates the meaning of the fields to the corresponding database.

In addition, the parsers are written so that the handling of separate fields can be easily overridden in a subclass. This makes it possible to extend them or take some local *specificities* into account (if you need to massage data that contains systematic errors, this proves *very* useful).

The following example assumes you have created a BibTeX-compatible database, as explained in the section 1.2.2 on page 5. It will then open a proper BibTeX file, and merge it into the current database. The list of imported references is returned as a result set.

```
from Pybl.io.Parsers.Semantic import BibTeX

parser = BibTeX.Reader()

rs = parser.parse(open('example.bib'), db)
```

1.4 Citation formatting

The *painful* part of writing formatting code is to take into account the missing fields without multiplying explicit checks that would quickly be boring. In addition, it is important to make it easy to factor out common operations, like formatting a list of authors, and reuse them in different contexts.

pybliographer provides a *domain specific language* that addresses these problems. However, it is not intended as a complete formatting language, so you cannot use it for instance to lay out your citations in a complete HTML web page (but this specific part is comparatively easy).

Back to practice. You can define some citation fragments like this:

```
from Pybl.io.Format import People, all, one

authors = People.lastFirst(all('author'))
title = one('title') | u'(no title)'
```

In this example, the `authors` variable is composed by taking all the values in the `author` field (`all('author')`), and pass them through the `lastFirst` transformation, which will format them as *Last Name, First Name*. The `Person` module contains other formatting variants for person names.

The `title` variable is built by taking the first value of the `title` field (via the `one` operator), and in case it does not exist, by using the string *no title* instead. This `|` alternative operator can be used everywhere a definition can be invalid.

You can then group these fields together, possibly adding some style information in the process:

```
from Pyblib.Format import B

citation = join(', ')[B[title], authors]
```

The `join` operator will take the parts between square braces and link them together with the text specified in parameter, a comma in that case. When one of the composing parts is not available, it is simply ignored, unless no part is available, in which case the whole expression is invalid (which can be trapped by using the `|` operator). In addition, the title is enclosed in a bold `B` tag.

Once the citation style is defined, it must be *compiled* on a specific database:

```
formatter = citation(db)
```

This operation checks that all the fields accessed are actually part of the schema. It also pre-computes certain information, so that the actual formatting of specific records can be a fast process.

Then, you can use the returned formatter and apply it to any number of records from the corresponding database:

```
cited = formatter(record)
```

You still need to get a definitive result, as you still need to select the output format for your citation. If you want it in HTML, you can do this last operation:

```
from Pyblib.Format import HTML

html = HTML.generate(cited)
```

Chapter 2

Extending *pybliographer*

TODO

2.1 Specializing a parser

TODO