

Francesc Alted

PyTables User's Guide



A hierarchical database for Python
Release 0.7.2

Alted, Francesc:

PyTables User's Guide

A hierarchical database for Python
Release 0.7.2

All rights reserved.
© 2002, 2003 Francesc Alted

Typeset by Francesc Alted
Day of print: 2003, September, 22th

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Copyright Notice and Statement for NCSA Hierarchical Data Format (HDF) Software Library and Utilities

NCSA HDF5 (Hierarchical Data Format 5) Software Library and Utilities Copyright 1998, 1999, 2000, 2001, 2002, 2003 by the Board of Trustees of the University of Illinois. All rights reserved.

See more information about the terms of this license at:
<http://hdf.ncsa.uiuc.edu/HDF5/doc/Copyright.html>

Copyright Notice and Statement for AURA *numarray* software library

Copyright (C) 2001 Association of Universities for Research in Astronomy (AURA)

THIS SOFTWARE IS PROVIDED BY AURA "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL AURA BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Contents

1	Introduction	1
1.1	Main Features	1
1.2	The Object Tree	2
2	Installation	7
2.1	Installation from sources	7
2.2	Binary installation (Windows)	9
2.2.1	Basic installation (i.e. with no LZO and UCL compression libraries support)	9
2.2.2	Advanced installation (i.e. with support for LZO and UCL compression libraries)	10
3	Some tutorials	11
3.1	Getting started	11
3.1.1	Importing <code>tables</code> objects	11
3.1.2	Declaring a Column Descriptor	12
3.1.3	Creating a <code>PyTables</code> file from scratch	12
3.1.4	Creating a new group	13
3.1.5	Creating a new table	13
3.1.6	Reading (and selecting) data in table	14
3.1.7	Creating new array objects	15
3.1.8	Closing the file and looking at its content	16
3.2	Browsing the <i>object tree</i> and more	17
3.2.1	Traversing the object tree	17
3.2.2	Setting and getting user attributes	19
3.2.3	Getting object metadata	22
3.2.4	Reading actual data from <code>Array</code> objects	24
3.2.5	Appending data to an existing table	25
3.2.6	And finally... how to remove rows from a table	26
3.3	Multidimensional table cells and automatic sanity checks	26
3.3.1	Shape checking	29
3.3.2	Field name checking	29
3.3.3	Data type checking	30
4	Library Reference	31
4.1	<code>tables</code> variables and functions	31
4.1.1	Global variables	31
4.1.2	Global functions	31
4.2	The <code>IsDescription</code> class	32
4.3	The <code>Col</code> class and its descendants	32
4.4	The <code>File</code> class	33
4.4.1	<code>File</code> instance variables	33
4.4.2	<code>File</code> methods	34
4.4.3	<code>File</code> special methods	36
4.5	The <code>Group</code> class	37

4.5.1	Group instance variables	37
4.5.2	Group methods	37
4.5.3	Group special methods	38
4.6	The Leaf class	39
4.6.1	Leaf instance variables	39
4.6.2	Leaf methods	39
4.7	The Table class	39
4.7.1	Table instance variables	39
4.7.2	Table methods	40
4.7.3	Table special methods	40
4.8	The Row class	42
4.8.1	Row methods	42
4.9	The Array class	42
4.9.1	Array instance variables	42
4.9.2	Array methods	42
4.10	The AttributeSet class	42
4.10.1	AttributeSet instance variables	43
4.10.2	AttributeSet methods	43
5	Optimization tips	45
5.1	Taking advantage of Psycho	45
5.2	Compression issues	46
5.3	Informing PyTables about expected number of rows in tables	49
5.4	Selecting an User Entry Point (UEP) in your tree	50
A	Supported data types in tables	51

*La sabiduría no vale la pena si no es
posible servirse de ella para inventar una
nueva manera de preparar los garbanzos.*

—Un sabio catalán
in "Cien años de soledad"
Gabriel García Márquez

Chapter 1

Introduction

The goal of PyTables is to enable the end user to manipulate easily scientific data **tables** and array objects in a hierarchical structure. The foundation of the underlying hierarchical data organization is the excellent HDF5 library (see NCSA).

It is important to remark that this package is not intended to serve as a complete wrapper for the entire HDF5 API, but to provide a flexible, *very Pythonic* tool to deal with (arbitrary) large amounts of data (typically bigger than available memory) in tables and arrays organized in a hierarchical, persistent disk storage.

A table is defined as a collection of records whose values are stored in *fixed-length* fields. All records have the same structure and all values in each field have the same *data type*. The terms *fixed-length* and strict *data types* seems to be quite a strange requirement for an interpreted language like Python, but they serve a useful function if the goal is to save very large quantities of data (such as is generated by many scientific applications, for example) in an efficient manner that reduces demand on CPU time and I/O.

In order to emulate records (that will be mapped to C structs in HDF5) in Python PyTables implements a special *metaclass* object in order to easily define all its fields and other properties. PyTables also provides a powerful interface to mine data in table. Records in tables are also known, in the HDF5 naming scheme, as *compound* data types.

For example, you can define arbitrary tables in Python simply by declaring a class with the name field and types information, like in:

```
class Particle(IsDescription):
    name      = StringCol(16)      # 16-character String
    idnumber  = Int64Col()         # Signed 64-bit integer
    ADCcount  = UInt16Col()        # Unsigned short integer
    TDCcount  = UInt8Col()         # unsigned byte
    grid_i    = Int32Col()         # integer
    grid_j    = IntCol()           # integer (equivalent to Int32Col)
    pressure  = Float32Col(shape=(2,3)) # 2-D float array (single-precision)
    energy    = FloatCol(shape=(2,3,4)) # 3-D float array (double-precision)
```

then, you have to pass this class to the table constructor, fill its rows with your values, and save (arbitrary large) collections of them in a file for persistent storage. After that, this data can be retrieved and post-processed quite easily with PyTables or even with another HDF5 application (in C, Fortran, Java or whatever language that provides an interface to HDF5).

Next section describes the most interesting capabilities of PyTables.

1.1 Main Features

PyTables take advantage of the powerful object orientation and introspection capabilities offered by Python to bring the next features to the user:

- *Support of table entities*: Allows working with a large number of records, i.e. that don't fit in memory.

- *Appendable tables:* It supports adding records to already created tables. This can be done without copying the dataset or redefining its structure, even between different Python sessions.
- *Multidimensional table cells:* You can declare a column to be formed by general array cells, in addition to only scalars, as the majority of relational databases do.
- *Support of arrays:* Numeric (see Ascher *et al.*) and numarray (see Greenfield *et al.*) arrays are a very useful complement of tables to keep homogeneous table slices (like selections of table columns).
- *Supports a hierarchical data model:* That way, you can structure very clearly all your data. PyTables builds up an *object tree* in memory that replicates the underlying file data structure. Access to the file objects is achieved by walking throughout this object tree, and manipulating it.
- *Support of files bigger than 2 GB:* The underlying HDF5 library already can do that (if your platform supports the C long long integer, or, on Windows, `__int64`), and PyTables automatically inherits this capability.
- *Can read generic HDF5 files:* PyTables can access to objects in generic HDF5 files provided they contain any combination of groups, compound type datasets (that will be mapped to Table objects) or homogeneous datasets (that will be mapped to Array objects). However, as these kind of data is the most common to be saved HDF5 format, PyTables can probably access most of the HDF5 files out there.
- *Data compression:* It supports data compression (through the use of the `zlib`, `LZO` and `UCL` libraries) out of the box. This become important when you have repetitive data patterns and don't want to loose your time searching for an optimized way to save them (i.e. it saves you data organization analysis time).
- *High performance I/O:* On modern systems, and for large amounts of data, tables and array objects can be read and written at a speed only limited by the performance of the underlying I/O subsystem. Moreover, if your data is compressible, even faster than that!
- *Architecture-independent:* PyTables has been carefully coded (as HDF5 itself) with little-endian/big-endian byte orderings issues in mind. So, in principle, you can write a file in a big-endian machine (like a Sparc or MIPS) and read it in other little-endian (like Intel or Alpha) without problems.

1.2 The Object Tree

The hierarchical model of the underlying HDF5 library allows PyTables to manage tables and arrays in a tree-like structure. In order to achieve this, an *object tree* entity is *dynamically* created imitating the HDF5 structure on disk. That way, the access to the HDF5 objects is made by walking throughout this object tree, and, by looking at their *metadata* nodes, you can get a nice picture of what kind data is kept there.

The different nodes in the object tree are instances of PyTables classes. There are several types of those classes, but the most important ones are the `Group` and the `Leaf`. `Group` instances (that we will be calling *groups* from now on) are a grouping structure containing instances of zero or more groups or leaves, together with supplementary metadata. `Leaf` instances (that will be called *leaves*) are containers for actual data and cannot contain further groups or leaves. The `Table` and `Array` classes are descendants of `Leaf`, and inherits all its properties.

Working with groups and leaves is similar in many ways to working with directories and files, respectively, in a Unix filesystem. As with Unix directories and files, objects in the object tree are often described by giving their full (or absolute) path names. In PyTables this full path can be specified either as string (like in `'/subgroup2/table3'`) or as a complete object path written in a certain way known as *natural name* schema (like in `file.root.subgroup2.table3`).

The support for *natural naming* is a key aspect of PyTables and means that the names of instance variables of the node objects are the same as the names of the element's children¹. This is very *Pythonic* and comfortable in many cases, as you can check in the tutorial section 3.1.6.

¹ I have got this simple but powerful idea from the excellent `Objectify` module by David Mertz (see Mertz)

You should also note that not all the data present on file is loaded in the object tree, but only the *metadata* (i.e. special data that describes the structure of the actual data). The actual data is not read until you ask for it (by calling a method on a particular node). By making use of the object tree (the metadata) you can get information on the objects on disk such as table names, title, name columns, data types in columns, the number of rows, or, in the case of arrays, the shape, the typecode, and so on. You can also traverse the tree in order to search for something and when you find the data you are interested in you can read it and process it. In some sense, you can think of PyTables as a tool that provide the same introspection capabilities of Python objects, but applied to the persistent storage of large amounts of data.

To better understand the dynamic nature of this object tree entity, let's start by a first example and try to realize what kind of object tree the next script (you can find it in `examples/objecttree.py`) would create:

```
from tables import *

class Particle(IsDescription):
    identity = StringCol(length=22, dflt=" ", pos = 0) # character String
    idnumber = Int16Col(1, pos = 1) # short integer
    speed    = Float32Col(1, pos = 1) # single-precision

# Open a file in "w"rite mode
fileh = openFile("objecttree.h5", mode = "w")
# Get the HDF5 root group
root = fileh.root

# Create the groups:
group1 = fileh.createGroup(root, "group1")
group2 = fileh.createGroup(root, "group2")

# Now, create a table in "group0" group
array1 = fileh.createArray(root, "array1", ["string", "array"], "String array")
# Create 2 new tables in group1
table1 = fileh.createTable(group1, "table1", Particle)
table2 = fileh.createTable("/group2", "table2", Particle)
# Create the last table in group2
array2 = fileh.createArray("/group1", "array2", [1,2,3,4])

# Now, fill the tables:
for table in (table1, table2):
    # Get the record object associated with the table:
    row = table.row
    # Fill the table with 10 records
    for i in xrange(10):
        # First, assign the values to the Particle record
        row['identity'] = 'This is particle: %2d' % (i)
        row['idnumber'] = i
        row['speed'] = i * 2.
        # This injects the Record values
        row.append()

    # Flush the table buffers
    table.flush()

# Finally, close the file (this also will flush all the remaining buffers!)
fileh.close()
```

This small program creates a simple HDF5 file, called `objecttree.h5`, with the structure that appears



Figure 1.1: An HDF5 example with 2 subgroups, 2 tables and 1 array.

in figure 1.1. During creation time, metadata in the object tree is updated in memory while the actual data is being saved on disk and when you close the file the object tree becomes unavailable. But, when you will open again this file the object tree will be re-constructed in memory from the metadata existent on disk, so that you can work with it exactly in the same way than during the original creation process.

In figure 1.2 you can see an example of the object tree created by reading the above `objecttree.h5` file (in fact, such an object is always created when reading any supported generic HDF5 file). If you are going to become a PyTables user, take your time to understand it². That will also make you more proactive by avoiding programming mistakes.

² Bear in mind, however, that this diagram is **not** a standard UML class diagram; it is rather meant to show the connections between the PyTables objects and some of its most important attributes and methods.

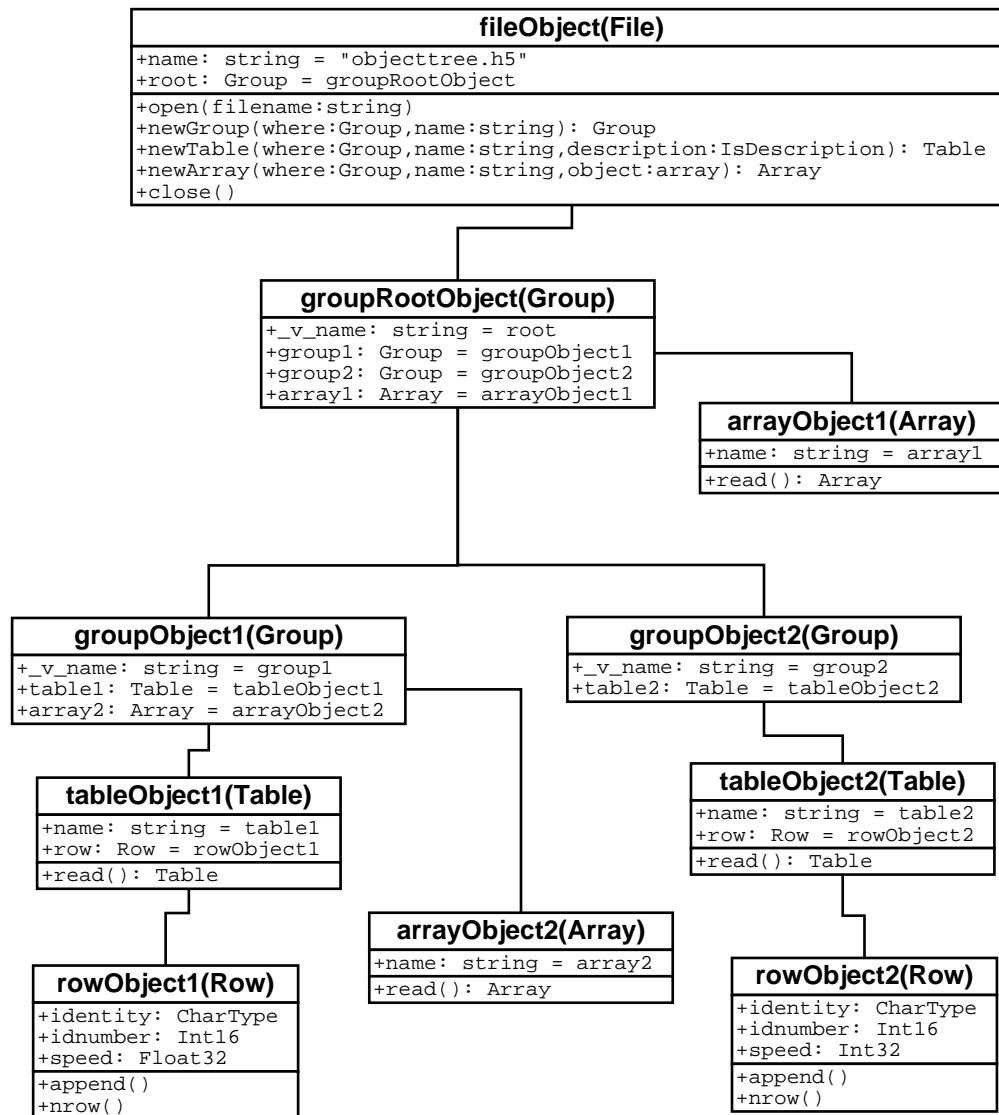


Figure 1.2: An object tree example in PyTables.

Chapter 2

Installation

*T'adones, company,
no volen arguments,
usen la força,
t'adones, amic.
T'adones, company,
que hem de sortir al carrer
junts, molts, com més millor,
si no volem perdre-ho tot,
t'adones, amic.*

—Raimon in the song "T'adones, amic"

The Python `Distutils` are used to build and install `PyTables`, so it is fairly simple to get things ready to go. Also, a binary distribution is available for Windows (see section 2.2).

2.1 Installation from sources

These instructions are both for Unix/Linux and Windows systems. If you are using Windows, it is assumed that you are using a recent version of MS Visual C++ (≥ 6.0) compiler. A GCC compiler is assumed for Unix, but other compilers should work as well.

Extensions in `PyTables` has been made using Pyrex (see Ewing) and C language. You can rebuild everything from scratch if you got Pyrex installed, but this is not necessary, as the Pyrex compiled source is included in the distribution. In order to do that, merely replace `setup.py` script in these instructions by `setup-pyrex.py`.

To compile `PyTables` you will need a recent version of HDF5 (C flavor) library and `numarray` (see Greenfield *et al.*) package. Although you won't need Numerical Python (see Ascher *et al.*) in order to compile `PyTables`, it is supported; you will only need a reasonably recent version of it ($\geq 21.x$). `PyTables` has been successfully tested with Numeric 21.3, 22.0 and 23.0. If you already have Numeric installed, the test driver module will detect it and will run the tests for Numeric automatically.

1. First, make sure that you have HDF5 1.6.x and `numarray` 0.6 or higher installed (I'm using HDF5 1.6.0 and `numarray` 0.7 currently). If don't, you can find them at <http://hdf.ncsa.uiuc.edu/HDF5> and <http://www.pfdubois.com/numpy>. Compile/install them.

Optionally, consider to install the excellent LZO and UCL compression libraries (see Oberhumer and section 5.2).

Unix `setup.py` will detect HDF5, LZO or UCL libraries and include files under `/usr` or `/usr/local`; this will catch installations from RPMs, DEBs and most hand installations under Unix. If `setup.py` can't find your `libhdf5` (or any other library you may wish) or if you have several versions installed and want to select one of them, then you can give it a hint either in the environment (using the `HDF5_DIR` environment variable or `LZO_DIR` and `UCL_DIR` for the optional libraries) or on the command line by specifying the directory containing the include and lib directory. For example:

```
--hdf5=/stuff/hdf5-1.6.0
--ucl=/stuff/ucl-1.0.1
```

If your HDF5 library was built as shared library, and if this shared library is not in the runtime load path, then you can specify the additional linker flags needed to find the shared library on the command line as well. For example:

```
--lflags="-Xlinker -rpath -Xlinker /stuff/hdf5-1.6.0/lib"
```

or perhaps just

```
--lflags="-R /stuff/hdf5-1.6.0/lib"
```

Check your compiler and linker documentation for correct syntax.

It is also possible to specify linking against different libraries with the `--libs` switch:

```
--libs="-lhdf5-1.6.5"
--libs="-lhdf5-1.6.5 -lnsl"
```

Windows If you are using Windows, and don't want to compile the libraries by hand, there are available binary packages for them. You can fetch the HDF5 and numarray binaries from their homes. For the HDF5 you must manually copy the `hdf5dll.dll` **and** `szipdll.dll` files to a directory in your `PATH` environment variable (for example `C:\WINDOWS\SYSTEM`).

Besides, you may want to fetch the LZO and UCL binaries from:

<http://gnuwin32.sourceforge.net/>. Normally, you will only need to fetch and install the `<package>-<version>-bin.zip` file, although in some cases the headers are in `<package>-<version>-lib.zip` file.

As before, you will need to copy manually the `lzo.dll` or `ucl.dll` to a directory in the `PATH`, so that they can be found by PyTables extensions.

`setup.py` will need that you inform it about where the library *stubs* (`.lib`) are installed as well as the *header* files (`.h`). To tell `setup.py` where the stubs and headers are, set the next environment variables:

HDF5_DIR Points to the HDF5 main directory (where the `include/` and `dll/` directories hangs).

Mandatory.

LZO_DIR Points to the LZO main directory (where the `include/` and `lib/` directories hangs).

Optional.

UCL_DIR Points to the UCL main directory (where the `include/` and `lib/` directories hangs).

Optional.

For example:

```
set HDF5_DIR=c:\stuff\5-160-win2k\c\release
set LZO_DIR=c:\stuff\lzo-1-07
set UCL_DIR=c:\stuff\ucl-1-01
```

Or you can pass this info to `setup.py` within the command line by specifying the directory containing the `include` and `lib` directory. For example:

```
--hdf5=c:\stuff\5-160-win2k\c\release --lzo=c:\stuff\lzo-1-07
--ucl=c:\stuff\ucl-1-01
```

2. From the main PyTables distribution directory run this command, (plus any extra flags needed as discussed above):

```
python setup.py build_ext --inplace
```

depending on the compiler flags used when compiling your Python executable, there may appear lots of warnings. Don't worry, almost all of them are caused by variables declared but never used. That's normal in Pyrex extensions.

3. To run the test suite change into the test directory and run this command:

Unix In the shell `sh` and its variants:

```
PYTHONPATH=..
export PYTHONPATH
python test_all.py
```

Windows Open a DOS terminal and write:

```
set PYTHONPATH=..
python test_all.py
```

If you would like to see some verbose output from the tests simply add the flag `-v` and/or the word `verbose` to the command line. You can also run just the tests in a particular test module. For example:

```
python test_types.py -v
```

If there is some test that do not pass, please, run the failing test module with all verbosity enabled (flags `-v verbose`), and send back the output to developers.

If you run into problems because Python can't load the HDF5, or any other shared library:

Unix Try to set the `LD_LIBRARY_PATH` environment variable to point to the directory where the libraries are.

Windows Put the DLL libraries (`hdf5dll.dll` and, optionally, `lzo.dll` and `ucl.dll`) on a directory listed on your `PATH` environment variable. The `setup.py` should already warned you about that.

4. To install the entire PyTables Python package, change back to the root distribution directory and run this command as the root user (remember to add any extra flags needed):

```
python setup.py install
```

That's it!. Now, proceed with the next chapter to see how to use PyTables.

2.2 Binary installation (Windows)

This section is only intended for Windows platforms. If you have Unix, or want to compile PyTables for Windows, jump to the section 2.1.

2.2.1 Basic installation (i.e. with no LZO and UCL compression libraries support)

First, make sure that you have HDF5 1.6.x or higher and numarray 0.6 or higher installed (I'm using HDF5 1.6.0 and numarray 0.6 currently). If don't, you can find them at <http://hdf.ncsa.uiuc.edu/HDF5> and <http://sourceforge.net/projects/numpy/>. Download the binary packages and install them. For the HDF5 it should be enough to manually copying the `hdf5dll.dll` **and** `szipdll.dll` files to a directory in your `PATH` environment variable (for example `C:\WINDOWS\SYSTEM`). Also, you should know that HDF5 **needs** the **zlib** compression library to work, so you will have to download and install it as well; a pre-compiled version is available at <ftp://ftp.ncsa.uiuc.edu/HDF/gzip/windows/>.

Caveat: When downloading the binary file for HDF5 libraries, select one compiled with the MSVC 6.0, like for example, the file `5-160-win2k.zip`, no matter if you are using Win2k or WinXP (it should work fine on both). The file `5-160-winxp-net.zip` was compiled with the MSVC 7.0 (aka ".NET") and **does not** work well with the PyTables binary (which has been generated with MSVC 6.0). You have been warned!.

Download the `tables-<version>.win32-py<version>.exe` file and execute it. You are done!.

You can (*you should*) test your installation by unpacking the source tarball, changing to the `test/` subdirectory and execute the `test_all.py` script. If all the tests passes (maybe with a couple of warnings, related with the possible unavailability of LZO and UCL libs, but that's ok for the binary version) you already have a working, well tested, copy of PyTables installed!. If don't, please, execute the `test_all.py -v verbose` and return the output to me.

2.2.2 Advanced installation (i.e. with support for LZO and UCL compression libraries)

If you want support for LZO and UCL libraries (see section 5.2 for hints about what they are useful for), install the LZO and UCL DLL libraries by fetching the binaries from:<http://gnuwin32.sourceforge.net/>. Normally, you will only need to fetch and install the <package>-<version>-bin.zip file.

Finally, follow the same steps than in section 2.2.1, but fetch and install the `tables-<version>-LU.win32-py<version>.exe` self-installer in place of `tables-<version>.win32-py<version>.exe`.

That's it!. Now, proceed with the next chapter to see how to use PyTables.

*Tout le malheur des hommes vient d'une
seule chose, qui est de ne savoir pas
demeurer en repos, dans une chambre.*

—Blaise Pascal

Chapter 3

Some tutorials

This chapter begins with a series of simple, yet comprehensive sections written in a tutorial style that will let you understand the main features that PyTables provide. If during the trip you want more information on some specific instance variable, global function or method, look at the doc strings or go to the library reference in chapter 4. However, if you are reading this in PDF or HTML formats, there should be an hyperlink to its reference near each newly introduced entity.

Please, note that throughout this document the terms *column* and *field* will be used interchangeably with the same meaning, and the same goes for the terms *row* and *record*.

3.1 Getting started

In this section, we will see how to define our own records from Python and save collections of them (i.e. a **table**) on a file. Then, we will select some data in the table using Python cuts, creating `numarray` arrays to keep this selection as separate objects in the tree.

In `examples/tutorial1-1.py` you will find the working version of all the code in this section. Nonetheless, this tutorial series has been written to allow you reproduce it in a Python interactive console. You are encouraged to take advantage of that by doing parallel testing and inspecting the created objects (variables, docs, children objects, etc.) during the voyage!.

3.1.1 Importing tables objects

Before doing anything you need to import the public objects in the `tables` package. You normally do that by issuing:

```
>>> import tables
>>>
```

This is the recommended way to import `tables` if you don't want to pollute too much your namespace. However, PyTables has a very reduced set of first-level primitives, so you may consider to use this alternative:

```
>>> from tables import *
>>>
```

that will export in your caller application namespace the next objects: `openFile`, `isHDF5`, `isPyTablesFile` and `IsDescription`. These are a rather small number of objects, and for convenience, we will use this last way to access them.

If you are going to deal with `numarray` or `Numeric` arrays (and normally, you will) you also need to import some objects from it. You can do that in the normal way. So, to access to PyTables functionality normally you should start your programs with:


```
>>> import tables          # but in this tutorial we use "from tables import *"
>>> from numarray import *  # or "from Numeric import *"
>>>
```

3.1.2 Declaring a Column Descriptor

Now, imagine that we have a particle detector and we want to create a table object in order to save data that comes from it. You need first to define that table, how many columns it have, which kind of object is each element on the columns, and so on.

Our detector has a TDC (Time to Digital Converter) counter with a dynamic range of 8 bits and an ADC (Analogic to Digital Converter) with a range of 16 bits. For these values, we will define 2 fields in our record object called `TDCcount` and `ADCcount`. We also want to save the grid position in which the particle has been detected and we will add two new fields called `grid_i` and `grid_j`. Our instrumentation also can obtain the pressure and energy of this particle that we want to add in the same way. The resolution of pressure-gauge allows us to use simple-precision float which will be enough to save pressure information, while energy would need a double-precision float. Finally, to track this particle we want to assign it a name to inform about the kind of the particle and a number identifier unique for each particle. So we will add a couple of fields: `name` will be the a string of up-to 16 characters and because we want to deal with a really huge number of particles, `idnumber` will be an integer of 64 bits.

With all of that, we can declare a new `Particle` class that will keep all this info:

```
>>> class Particle(IsDescription):
...     name      = StringCol(16)      # 16-character String
...     idnumber  = Int64Col()         # Signed 64-bit integer
...     ADCcount  = UInt16Col()        # Unsigned short integer
...     TDCcount  = UInt8Col()         # unsigned byte
...     grid_i    = Int32Col()         # integer
...     grid_j    = IntCol()           # integer (equivalent to Int32Col)
...     pressure  = Float32Col()       # float (single-precision)
...     energy    = FloatCol()         # double (double-precision)
...
>>>
```

This definition class is quite auto-explanatory. Basically, you have to declare a class variable for each field you need, and as its value we assign a subclass instance of the `Col` class, that describes the kind of column (the data type, the length, the shape, ...). See section 4.3 for a complete description of these subclasses. See also appendix A for a list of data types supported in `Col` constructors.

From now on, we can use `Particle` instances as a descriptor for our detector data table. We will see how to pass this object to the `Table` constructor. But first, we must create a file where all the actual data pushed into `Table` will be saved.

3.1.3 Creating a PyTables file from scratch

To create a PyTables file use the first-level `openFile` (see 4.1.2) function:

```
>>> h5file = openFile("tutorial1.h5", mode = "w", title = "Test file")
```

This `openFile` (see 4.1.2) is one of the objects imported by the `"from tables import *"`, do you remember?. Here, we are telling that we want to create a new file called `"tutorial1.h5"` in `"w"`rite mode and with an descriptive title string (`"Test file"`). This function tries to open the file, and if successful, returns a `File` (see 4.4) instance which hosts the root of the object tree on its `root` attribute.

3.1.4 Creating a new group

Now, to better organize our data, we will create a group hanging from the root called *detector*. We will use this group to save our particle data there.

```
>>> group = h5file.createGroup("/", 'detector', 'Detector information')
>>>
```

Here, we have taken the `File` instance `h5file` and invoked its `createGroup` method (see 4.4.2), telling that we want to create a new group called *detector* hanging from `"/"`, which is other way to refer to the `h5file.root` object we mentioned before. This will create a new `Group` (see 4.5) instance that will be assigned to the `group` variable.

3.1.5 Creating a new table

Let's now create the `Table` (see 4.7) object hanging from the new created group. We do that by calling the `createTable` (see 4.4.2) method from the `h5file` object:

```
>>> table = h5file.createTable(group, 'readout', Particle, "Readout example")
>>>
```

Look at how we asked to create the `Table` instance hanging from `group`, with name *"readout"*. We have passed `Particle`, the class that we have declared before, as the *description* parameter and finally we have used *"Readout example"* as a `Table` title. With all this information, a new `Table` instance is created and assigned to *table* variable.

If you are getting curious how the object tree looks like at this moment, simply print the name of the `File` instance, *h5file*, and look at their output:

```
>>> print h5file
Filename: 'tutorial1.h5' Title: 'Test file' Last modif.: 'Sun Jul 27 14:00:13 2003'
/ (Group) 'Test file'
/detector (Group) 'Detector information'
/detector/readout (Table(0,)) 'Readout example'

>>>
```

As you can see, a dump of the object tree has been shown and it's very easy to visualize the `Group` and `Table` objects we have just created. If you want more information, just type the name of the `File` instance:

```
>>> h5file
>>> h5file
Filename: 'tutorial1.h5' Title: 'Test file' Last modif.: 'Sun Jul 27 14:00:13 2003'
/ (Group) 'Test file'
/detector (Group) 'Detector information'
/detector/readout (Table(0,)) 'Readout example'

>>> h5file
File(filename='tutorial1.h5', title='Test file', mode='w', trMap={}, rootUEP='/')
/ (Group) 'Test file'
/detector (Group) 'Detector information'
/detector/readout (Table(0,)) 'Readout example'
  description := {
    "ADCcount": Col('UInt16', shape=1, itemsize=2, dflt=0),
    "TDCcount": Col('UInt8', shape=1, itemsize=1, dflt=0),
    "energy": Col('Float64', shape=1, itemsize=8, dflt=0.0),
```

```
"grid_i": Col('Int32', shape=1, itemsize=4, dflt=0),
"grid_j": Col('Int32', shape=1, itemsize=4, dflt=0),
"idnumber": Col('Int64', shape=1, itemsize=8, dflt=0),
"name": Col('CharType', shape=1, itemsize=16, dflt=None),
"pressure": Col('Float32', shape=1, itemsize=4, dflt=0.0) }
byteorder := little

>>>
```

where more detailed info is printed on each object on the tree. Pay attention on how `Particle`, our table descriptor class, is printed as part of the *readout* table description information. In general, you can obtain lot of information on the objects and its children by just printing them. That introspection capability is very meaningful, so I recommend you to use it extensively.

Now, time to fill this table with some values. But first, we are going to get a pointer to the `Row` instance of this table instance:

```
>>> particle = table.row
>>>
```

The `row` attribute of `table` points to the `Row` (see 4.8) instance that will be used to input data rows into the table. We achieve this by just assigning it the values for each row as if it was a dictionary (although it is actually an *extension class*) and using the column names as keys.

Look at how the filling process works like:

```
>>> particle = table.row
>>> for i in xrange(10):
...     particle['name'] = 'Particle: %6d' % (i)
...     particle['TDCcount'] = i % 256
...     particle['ADCcount'] = (i * 256) % (1 << 16)
...     particle['grid_i'] = i
...     particle['grid_j'] = 10 - i
...     particle['pressure'] = float(i*i)
...     particle['energy'] = float(particle['pressure'] ** 4)
...     particle['idnumber'] = i * (2 ** 34)
...     particle.append()
...
>>>
```

This code should be easy to understand. The lines inside the loop just assign values to the different columns in the `particle` `Row` instance (see 4.8) and then a call to its `append()` method is made to put this information in the `table` I/O buffer.

After we have pushed all our data, we should flush the I/O buffer for the table if we want to consolidate all this data on disk. We can achieve that by calling the `table.flush()` method.

```
>>> table.flush()
>>>
```

3.1.6 Reading (and selecting) data in table

Ok. We have now our data on disk but to this data be useful we need to access it and select some values we are interested in and located at some specific columns. That's is easy to do:

```
>>> table = h5file.root.detector.readout
>>> pressure = [ x['pressure'] for x in table.iterrows()
...             if x['TDCcount']>3 and 20<=x['pressure']<50 ]
>>> pressure
[25.0, 36.0, 49.0]
>>>
```

The first line is only to declare a convenient shortcut to the *readout* table which is a bit deeper on the object tree. As you can see, we have used the **natural naming** schema to access it. We could also have used the `h5file.getNode()` method instead, and we will certainly do that later on.

You will recognize the last two lines to be a Python list comprehension. It loops over rows in *table* as they are provided by `table.iterrows()` iterator (see 4.7.2) that returns values until data in table is exhausted. These rows are filtered using the expression `x['TDCcount'] > 3 and x['pressure'] < 50`, and the pressure field for satisfying records is selected to form the final list that is assigned to *pressure* variable.

We could indeed have used a normal `for` loop to do that, but I find comprehension syntax to be more compact and elegant.

Let's select the names for the same set of cuts:

```
>>> names=[ x['name'] for x in table if x['TDCcount']>3 and 20<=x['pressure']<50 ]
>>> names
['Particle:      5', 'Particle:      6', 'Particle:      7']
>>>
```

Note how we have omitted the `iterrows()` call in the list comprehension. This is because the *Table* class has an implementation of the special method called `__iter__()`, so that it implements the iterator protocol over all the rows in the table. In fact, `iterrows()` internally calls this special `__iter__()` method. This way to access all the rows in a table turns out to be very convenient, specially for interactive use.

Ok. that's enough for selections. Next section will show you how to save these selections on file.

3.1.7 Creating new array objects

In order to separate the selected data from the detector data, we will create a new group, called *columns* hanging from the root group:

```
>>> gcolumns = h5file.createGroup(h5file.root, "columns", "Pressure and Name")
>>>
```

Note that this time we have specified the first parameter in a *natural naming* fashion (`h5file.root`) instead of using an absolute path string (`"/"`).

Now, create one *Array* object:

```
>>> h5file.createArray(gcolumns, 'pressure', array(pressure),
...                  "Pressure column selection")
/cOLUMNS/pressure (Array(3,)) 'Pressure column selection'
  type = Float64
  itemsize = 8
  flavor = 'NumArray'
  byteorder = 'little'
>>>
```

We already know the first two parameters of the `createArray` (see 4.4.2) methods (these are the same as the firsts in `createTable`): they are the parent group *where* *Array* will be created and the *Array* instance

name. You can figure out that the fourth parameter is the *title*. And in the third position we have the *object* we want to save on disk. In this case, it is a Numeric array that is built from the selection lists we created before.

Now, we are going to save the other selection. In this case it's a list of strings, and we want to save this object as is, with no further conversion. Look at how this can be done:

```
>>> h5file.createArray(gcolumns, 'name', names, "Name column selection")
/cOLUMNS/name Array(4,) 'Name column selection'
  type = 'CharType'
  itemsizes = 16
  flavor = 'List'
  byteorder = 'little'
>>>
```

You see, `createArray()` accepts *names* (which is a regular Python list) as *object* parameter. Actually, it accepts a variety of other regular objects (see 4.4.2). We will check that we can retrieve exactly the same object from disk later on.

Note that in this examples, `createArray` method returns an `Array` instance that is not assigned to any variable. Don't worry, this was intentional because I wanted to show you the kind of object we have created by showing its representation. Indeed, the `Array` objects has been attached to the object tree and saved on disk, as you can see if you print the complete object tree:

```
>>> print h5file
Filename: 'tutorial1.h5' Title: 'Test file' Last modif.: 'Sun Jul 27 14:00:13 2003'
/ (Group) 'Test file'
/cOLUMNS (Group) 'Pressure and Name'
/cOLUMNS/name (Array(3,)) 'Name column selection'
/cOLUMNS/pressure (Array(3,)) 'Pressure column selection'
/detector (Group) 'Detector information'
/detector/readout (Table(10,)) 'Readout example'

>>>
```

3.1.8 Closing the file and looking at its content

To finish this first tutorial, we use the `close` method of the `h5file File` instance to close the file before exiting Python:

```
>>> h5file.close()
>>> ^D
```

With all that, you have created your first PyTables file with a table and two arrays. That was easy, admit it. Now, you can have a look at it with some generic HDF5 tool, like `h5dump` or `h5ls`. Here is the result of passing to `h5ls` the `tutorial1.h5` file:

```
$ h5ls -rd tutorial1.h5
/cOLUMNS                               Group
/cOLUMNS/name                         Dataset {3}
  Data:
    (0) "Particle:      5", "Particle:      6", "Particle:      7"
/cOLUMNS/pressure                     Dataset {3}
  Data:
    (0) 25, 36, 49
/detector                               Group
```

```

/detector/readout      Dataset {10/Inf}
  Data:
    (0) {0, 0, 0, 0, 10, 0, "Particle:      0", 0},
    (1) {256, 1, 1, 1, 9, 17179869184, "Particle:      1", 1},
    (2) {512, 2, 256, 2, 8, 34359738368, "Particle:      2", 4},
    (3) {768, 3, 6561, 3, 7, 51539607552, "Particle:      3", 9},
    (4) {1024, 4, 65536, 4, 6, 68719476736, "Particle:      4", 16},
    (5) {1280, 5, 390625, 5, 5, 85899345920, "Particle:      5", 25},
    (6) {1536, 6, 1679616, 6, 4, 103079215104, "Particle:      6", 36},
    (7) {1792, 7, 5764801, 7, 3, 120259084288, "Particle:      7", 49},
    (8) {2048, 8, 16777216, 8, 2, 137438953472, "Particle:      8", 64},
    (9) {2304, 9, 43046721, 9, 1, 154618822656, "Particle:      9", 81}

```

or, using the "dumpFile.py" PyTables utility (located in examples/ directory):

```

$ python dumpFile.py tutorial1.h5
Filename: 'tutorial1.h5' Title: 'Test file' Last modif.: 'Sun Jul 27 14:40:51 2003'
/ (Group) 'Test file'
/columns (Group) 'Pressure and Name'
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
/detector (Group) 'Detector information'
/detector/readout (Table(10,)) 'Readout example'

```

You can pass the `-v` or `-d` options to `dumpFile.py` if you want more verbosity. Try them out!

3.2 Browsing the *object tree* and more

In this section, we will learn how to browse the tree while retrieving meta-information about the actual data, and will finish by appending some rows to the existing table to show how table objects can be enlarged.

In *examples/tutorial1-2.py* you will find the working version of all the code in this section. As before, you are encouraged to use a python shell and inspect the object tree during the voyage.

3.2.1 Traversing the object tree

First of all, let's open the file we have recently created in last tutorial section, as we will take it as a basis for this section:

```
>>> h5file = openFile("tutorial1.h5", "a")
```

This time, we have opened the file in "a"ppend mode. We are using this mode because we want to add more information to the file.

PyTables, following the Python tradition, offers powerful introspection capabilities, i.e. you can easily ask information about any component of the object tree as well as traverse the tree searching for something.

To start with, you can get a first glance image of the object tree, by simply printing the existing File instance:

```

>>> print h5file
Filename: 'tutorial1.h5' Title: 'Test file' Last modif.: 'Sun Jul 27 14:40:51 2003'
/ (Group) 'Test file'
/columns (Group) 'Pressure and Name'
/columns/name (Array(3,)) 'Name column selection'

```

```
/columns/pressure (Array(3,)) 'Pressure column selection'
/detector (Group) 'Detector information'
/detector/readout (Table(10,)) 'Readout example'

>>>
```

That's right, it seems that all our objects are there. Now, let's make use of the File iterator to see how to list all the nodes in the object tree:

```
>>> for node in h5file:
...     print node
...
/ (Group) 'Test file'
/columns (Group) 'Pressure and Name'
/detector (Group) 'Detector information'
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
/detector/readout (Table(10,)) 'Readout example'

>>>
```

We can use the `walkGroups` method (see 4.4.2) of File class to list only the *groups* on tree:

```
>>> for group in h5file.walkGroups("/"):
...     print group
...
/ (Group) 'Test file'
/columns (Group) 'Pressure and Name'
/detector (Group) 'Detector information'

>>>
```

Note that `walkGroups()` actually returns an *iterator*, not a list of objects. Combining this iterator with the `listNodes()` method, we can do very powerful things. Let's see an example listing all the arrays in the tree:

```
>>> for group in h5file.walkGroups("/"):
...     for array in h5file.listNodes(group, classname = 'Array'):
...         print array
...
/columns/name Array(4,) 'Name column selection'
/columns/pressure Array(4,) 'Pressure column selection'
```

`listNodes()` (see 4.4.2) returns a list containing all the nodes hanging from a specific Group, and if *classname* keyword is specified, the method will filter all instances which are not descendants of it. We have specified it to solely return Array instances.

We can combine both calls by using the `__call__`(where, classname) special method of File (see 4.4.3), i.e.:

```
>>> for array in h5file("/", "Array"):
...     print array
...
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'

>>>
```

which is a nice shortcut for doing interactive work.

As a final example, we will list all the Leaf, i.e. Table and Array instances (see 4.6 for detailed information on Leaf class), in /detector group. Check that only one instance of Table class (i.e. readout) will be selected in this group (as it should be):

```
>>> for leaf in h5file.root.detector('Leaf'):
...     print leaf
...
/detector/readout (Table(10,)) 'Readout example'
>>>
```

where we have used a call to the Group.__call__(classname, recursive) special method (4.5.3), combined with a *natural naming* path specification.

Of course you can do more sophisticated node selections using these powerful methods, but first, we need to learn a bit about some important instance variables of PyTables objects.

3.2.2 Setting and getting user attributes

PyTables provides an easy and concise way to complement the meaning of your node objects on the tree by using the AttributeSet class (see section 4.10). You can access to this object through the standard attribute attrs in Leaf nodes and _v_attrs in Group nodes.

For example, let's imagine that we want to save the date indicating when the data in /detector/readout table has been acquired, as well as the temperature during the gathering process. That is easy:

```
>>> table = h5file.root.detector.readout
>>> table.attrs.gath_date = "Wed, 06/12/2003 18:33"
>>> table.attrs.temperature = 18.4
>>> table.attrs.temp_scale = "Celsius"
>>>
```

Now, set a somewhat more complex attribute in the /detector group:

```
>>> detector = h5file.root.detector
>>> detector._v_attrs.stuff = [5, (2.3, 4.5), "Integer and tuple"]
>>>
```

Note how the AttributeSet instance is accessed with _v_attrs because detector is a Group node. In general, you can save any standard Python data structure as an attribute node, but see section 4.10 for a more detailed explanation of how this are serialized on disk.

Now, getting the attributes is equally easy:

```
>>> table.attrs.gath_date
'Wed, 06/12/2003 18:33'
>>> table.attrs.temperature
18.399999999999999
>>> table.attrs.temp_scale
'Celsius'
>>> detector._v_attrs.stuff
[5, (2.2999999999999998, 4.5), 'Integer and tuple']
>>>
```

You can probably guess how to delete attributes:

```
>>> del table.attrs.gath_date
```


If you want to have a look at the current attribute set of `/detector/table`, you can print its representation (try also hitting the `TAB` key twice if you are on a Unix Python console with the `rlcompleter` module active):

```
>>> table.attrs
/detector/readout (AttributeSet), 14 attributes:
[CLASS := 'TABLE',
 FIELD_0_NAME := 'ADCcount',
 FIELD_1_NAME := 'TDCcount',
 FIELD_2_NAME := 'energy',
 FIELD_3_NAME := 'grid_i',
 FIELD_4_NAME := 'grid_j',
 FIELD_5_NAME := 'idnumber',
 FIELD_6_NAME := 'name',
 FIELD_7_NAME := 'pressure',
 NROWS := 10,
 TITLE := 'Readout example',
 VERSION := '2.0',
 tempScale := 'Celsius',
 temperature := 18.399999999999999]
>>>
```

You can get a list only the user or system attributes with the `_v_list()` method.

```
>>> print table.attrs._f_list("user")
['temp_scale', 'temperature']
>>> print table.attrs._f_list("sys")
['CLASS', 'FIELD_0_NAME', 'FIELD_1_NAME', 'FIELD_2_NAME', 'FIELD_3_NAME',
 'FIELD_4_NAME', 'FIELD_5_NAME', 'FIELD_6_NAME', 'FIELD_7_NAME', 'NROWS',
 'TITLE', 'VERSION']
>>>
```

And rename attributes:

```
>>> table.attrs._f_rename("temp_scale", "tempScale")
>>> print table.attrs._f_list()
['tempScale', 'temperature']
>>>
```

However, you can't set, delete or rename read-only attributes:

```
>>> table.attrs._f_rename("VERSION", "version")
Traceback (most recent call last):
  File ">stdin>", line 1, in ?
  File "/home/falted/PyTables/pytables-0.7/tables/AttributeSet.py", line 249, in _f_rename
    raise RuntimeError, \
RuntimeError: Read-only attribute ('VERSION') cannot be renamed
>>>
```

After your session, you can check that the `/detector/readout` attributes in disk looks like:

```
$ h5ls -vr tutorial1.h5/detector/readout
Opened "tutorial1.h5" with sec2 driver.
/detector/readout      Dataset {10/Inf}
```

```

Attribute: CLASS      scalar
  Type:      6-byte null-terminated ASCII string
  Data:      "TABLE"
Attribute: VERSION    scalar
  Type:      4-byte null-terminated ASCII string
  Data:      "2.0"
Attribute: TITLE      scalar
  Type:      16-byte null-terminated ASCII string
  Data:      "Readout example"
Attribute: FIELD_0_NAME scalar
  Type:      9-byte null-terminated ASCII string
  Data:      "ADCcount"
Attribute: FIELD_1_NAME scalar
  Type:      9-byte null-terminated ASCII string
  Data:      "TDCcount"
Attribute: FIELD_2_NAME scalar
  Type:      7-byte null-terminated ASCII string
  Data:      "energy"
Attribute: FIELD_3_NAME scalar
  Type:      7-byte null-terminated ASCII string
  Data:      "grid_i"
Attribute: FIELD_4_NAME scalar
  Type:      7-byte null-terminated ASCII string
  Data:      "grid_j"
Attribute: FIELD_5_NAME scalar
  Type:      9-byte null-terminated ASCII string
  Data:      "idnumber"
Attribute: FIELD_6_NAME scalar
  Type:      5-byte null-terminated ASCII string
  Data:      "name"
Attribute: FIELD_7_NAME scalar
  Type:      9-byte null-terminated ASCII string
  Data:      "pressure"
Attribute: tempScale scalar
  Type:      8-byte null-terminated ASCII string
  Data:      "Celsius"
Attribute: temperature {1}
  Type:      native double
  Data:      18.4
Attribute: NROWS      {1}
  Type:      native int
  Data:      10
Location: 0:1:0:1952
Links: 1
Modified: 2003-07-24 13:59:19 CEST
Chunks: {2048} 96256 bytes
Storage: 470 logical bytes, 96256 allocated bytes, 0.49% utilization
Type: struct {
    "ADCcount"      +0    native unsigned short
    "TDCcount"      +2    native unsigned char
    "energy"        +3    native double
    "grid_i"        +11   native int
    "grid_j"        +15   native int
    "idnumber"      +19   native long long

```

```
        "name"                +27    16-byte null-terminated ASCII string
        "pressure"            +43    native float
    } 47 bytes
```

As you can see, the use of attributes can be a good mechanism to add persistent (meta) information to your actual data. Be sure to use them extensively.

3.2.3 Getting object metadata

Each object in PyTables has *metadata* information about the actual data on the file. Normally this *meta-information* is accessible through the node instance variables. Let's have a look at some examples:

```
>>> print "Object:", table
Object: /detector/readout Table(10,) 'Readout example'
>>> print "Table name:", table.name
Table name: readout
>>> print "Table title:", table.title
Table title: Readout example
>>> print "Number of rows in table:", table.nrows
Number of rows in table: 10
>>> print "Table variable names with their type and shape:"
Table variable names with their type and shape:
>>> for name in table.colnames:
...     print name, ' := %s, %s' % (table.coltypes[name], table.colshapes[name])
...
ADCcount := UInt16, 1
TDCcount := UInt8, 1
energy := Float64, 1
grid_i := Int32, 1
grid_j := Int32, 1
idnumber := Int64, 1
name := CharType, 1
pressure := Float32, 1
>>>
```

Here, the name, title, nrows, colnames, coltypes and colshapes attributes (see 4.4.1 for a complete attribute list) of Table object give us quite a lot of information about actual table data.

In general, you can get up-to-the-minute information about the public objects in PyTables in a interactive way by printing its internal doc strings:

```
>>> print table.__doc__
Represent a table in the object tree.
```

It provides methods to create new tables or open existing ones, as well as to write/read data to/from table objects over the file. A method is also provided to iterate over the rows without loading the entire table or column in memory.

Data can be written or read both as Row() instances or as numarray (NumArray or RecArray) objects.

Methods:

Common to all leaves:

```
close()
flush()
getAttr(attrname)
rename(newname)
remove()
setAttr(attrname, attrvalue)
```

Specific of Table:

```
iterrows()
read([start] [, stop] [, step] [, field [, flavor]])
removeRows(start, stop)
```

Instance variables:

Common to all leaves:

```
name -- the leaf node name
hdf5name -- the HDF5 leaf node name
title -- the leaf title
shape -- the leaf shape
byteorder -- the byteorder of the leaf
```

Specific of Table:

```
description -- the metaobject describing this table
row -- a reference to the Row object associated with this table
nrows -- the number of rows in this table
rowsize -- the size, in bytes, of each row
colnames -- the field names for the table (list)
coltypes -- the type class for the table fields (dictionary)
colshapes -- the shapes for the table fields (dictionary)
```

>>>

This is very handy if you don't have this manual at hand. Try yourself with other objects docs, like for example:

```
>>> help(table.__class__)
>>> help(table.removeRows)
```

Now, print some metadata in */columns/pressure* Array object:

```
>>> pressureObject = h5file.getNode("/columns", "pressure")
>>> print "Info on the object:", repr(pressureObject)
Info on the object: /columns/pressure (Array(3,)) 'Pressure column selection'
  type = Float64
  itemsize = 8
  flavor = 'NumArray'
  byteorder = 'little'
>>> print "  shape: ==>", pressureObject.shape
  shape: ==> (3,)
>>> print "  title: ==>", pressureObject.title
  title: ==> Pressure column selection
>>> print "  type: ==>", pressureObject.type
```

```
type: ==> Float64
>>>
```

Observe how we have used the `getNode()` method of `File` class to access a node in the tree, instead of the natural naming method. Both are useful, and depending on the context you will prefer to use one or another. `getNode()` has the advantage that it can get a node from the pathname string (like in this example), and, besides, you can force a filter so that the node in that location has to be a *classname* instance. However, I consider natural naming to be more elegant and quicker to specify, specially if you are using the name completion capability present in interactive console. I suggest to give a try at this powerful combination of natural naming and completion capabilities present on most Python consoles. You will see how pleasant can be browsing the object tree (well, as long as this activity can be qualified in that way).

If you look at the `type` attribute of the `pressureObject`, you can certify that this is a "**Float64**" array, and that by looking at their `shape` attribute, it can deduced that the array on disk is unidimensional and has 4 elements. See 4.9.1 or the internal string docs for the complete `Array` attribute list.

3.2.4 Reading actual data from Array objects

Once you have found the desired `Array` and decided that you want to retrieve the actual data array from it, you should use the `read()` method of the `Array` object:

```
>>> pressureArray = pressureObject.read()
>>> pressureArray
array([ 25.,  36.,  49.])
>>> print "pressureArray is an object of type:", type(pressureArray)
pressureArray is an object of type: <class 'numarray.numarraycore.NumArray'>
>>> nameArray = h5file.root.columns.name.read()
>>> nameArray
['Particle:      5', 'Particle:      6', 'Particle:      7']
>>> print "nameArray is an object of type:", type(nameArray)
nameArray is an object of type: <type 'list'>
>>>
>>> print "Data on arrays nameArray and pressureArray:"
Data on arrays nameArray and pressureArray:
>>> for i in range(pressureObject.shape[0]):
...     print nameArray[i], "-->", pressureArray[i]
...
Particle:      5 --> 25.0
Particle:      6 --> 36.0
Particle:      7 --> 49.0
>>> pressureObject.name
'pressure'
>>>
```

You can verify as the `read()` method (see section 4.9.2) returns an authentic `numarray` object for the `pressureObject` instance by looking at the output of the `type()` call, while for the `nameObject` instance `read()` returns a native Python list (of strings). This is because the type of the object saved is kept as an HDF5 attribute (named `FLAVOR`) for these objects on disk. This attribute is then read as part of the `Array` metainformation and accessible through the `Array.attrs.FLAVOR` variable, enabling the read array to be converted into the original object. This provides a means to save a large variety of objects as arrays with the guarantee that you will be able to recover them in its original form afterwards. See section 4.4.2 for a complete list of supported objects for `Array`.

3.2.5 Appending data to an existing table

Now, let's have a look at how we can add records to an existing on-disk table. Let's use our well-known *readout* Table instance and let's append some new values to it:

```
>>> table = h5file.root.detector.readout
>>> particle = table.row
>>> for i in xrange(10, 15):
...     particle['name'] = 'Particle: %6d' % (i)
...     particle['TDCcount'] = i % 256
...     particle['ADCcount'] = (i * 256) % (1 << 16)
...     particle['grid_i'] = i
...     particle['grid_j'] = 10 - i
...     particle['pressure'] = float(i*i)
...     particle['energy'] = float(particle['pressure'] ** 4)
...     particle['idnumber'] = i * (2 ** 34)
...     particle.append()
...
>>> table.flush()
>>>
```

That works exactly in the same way than filling a new table. PyTables knows that this table is on disk, and when you add new records, they are appended to the end of the table¹.

If you look carefully at the code you will see that we have used the `table.row` attribute so as to access a table row and fill it up with the new values. Each time that its `append()` method is called, the actual row is committed to the output buffer and the row pointer is incremented to point to the next table record. When the buffer is full, the data is saved on disk, and the buffer is reused again for the next cycle.

Caveat emptor!: Do not forget to always call the `.flush()` method after a writing operation; else your tables will not be fully updated!

Let's have a look at some columns of the resulting table:

```
>>> for r in table.iterrows():
...     print "%-16s | %11.1f | %11.4g | %6d | %6d | %8d |" % \
...         (r['name'], r['pressure'], r['energy'], r['grid_i'], r['grid_j'],
...         r['TDCcount'])
...
...
Particle:      0 |      0.0 |      0 |      0 |      10 |      0 |
Particle:      1 |      1.0 |      1 |      1 |       9 |      1 |
Particle:      2 |      4.0 |    256 |      2 |       8 |      2 |
Particle:      3 |      9.0 |   6561 |      3 |       7 |      3 |
Particle:      4 |     16.0 | 6.554e+04 |      4 |       6 |      4 |
Particle:      5 |     25.0 | 3.906e+05 |      5 |       5 |      5 |
Particle:      6 |     36.0 | 1.68e+06 |      6 |       4 |      6 |
Particle:      7 |     49.0 | 5.765e+06 |      7 |       3 |      7 |
Particle:      8 |     64.0 | 1.678e+07 |      8 |       2 |      8 |
Particle:      9 |     81.0 | 4.305e+07 |      9 |       1 |      9 |
Particle:     10 |    100.0 | 1e+08 |     10 |       0 |     10 |
Particle:     11 |    121.0 | 2.144e+08 |     11 |      -1 |     11 |
Particle:     12 |    144.0 | 4.3e+08 |     12 |      -2 |     12 |
Particle:     13 |    169.0 | 8.157e+08 |     13 |      -3 |     13 |
Particle:     14 |    196.0 | 1.476e+09 |     14 |      -4 |     14 |
```

¹ Note that you can append not only scalar values to tables, but also fully multidimensional array objects.

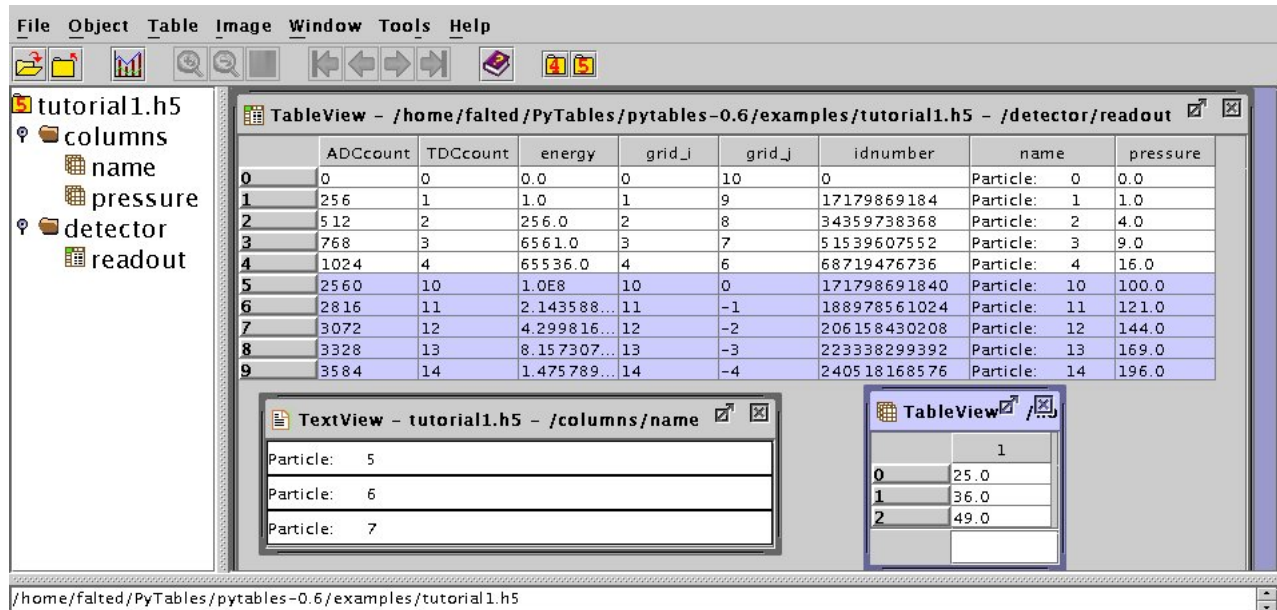


Figure 3.1: The final version of data file for tutorial 1, with a view of the data objects.

3.2.6 And finally... how to remove rows from a table

Let's start finishing this tutorial by deleting some rows from the table we have. Suppose that we want to delete the rows from 5th to 9th (inclusive). That's very easy to do:

```
>>> table.removeRows(5,10)
5
>>>
```

`removeRows(start, stop)` (see 4.7.2) deletes the rows in the range (start, stop). It returns the number of rows effectively removed.

We have reached the end of this first tutorial. But, ei!, do not forget to close the file after you finish all the work:

```
>>> h5file.close()
>>> ^D
$
```

In figure 3.1 you can see a graphical view of the PyTables file, with the datasets we have just created. And in figure 3.2 you can see the general properties of the table `/detector/readout`.

3.3 Multidimensional table cells and automatic sanity checks

Now, time for a more real life example (i.e. with errors in code). Here, we will create a couple of groups hanging directly from root called `Particles` and `Events`. Then, we will put 3 tables in each group; in `Particles` we will put tables based on `Particle` descriptor and in `Events`, tables based on `Event` descriptor.

After that, we will feed the tables with a number of records. Finally, we will read the recently created table `/Events/TEvent3` and select some values from it using a comprehension list.

Look at the next script (you can find it in `examples/tutorial2.py`). It seems to do all of that, but a couple of small bugs will be shown up. Note that this `Particle` class is not directly related with the one

General **Attributes**

Name: readout
 Path: /detector/
 Type: ncsa.hdf.object.h5.H5CompoundDS
 Object ID: 1952

Datatype

No. of Dimension(s): 1
 Dimension Size(s): 10
 Data Type: Compound

Name	Type	Array Size
ADCcount	16-bit unsigned integer	1
TDCcount	8-bit unsigned charac...	1
energy	64-bit floating-point	1
grid_i	32-bit integer	1
grid_j	32-bit integer	1
idnumber	64-bit integer	1
name	String, length=16	1
pressure	32-bit floating-point	1

Chunking: 2048
 Compression: NONE

Close

Figure 3.2: General properties of the /detector/readout table.

defined in last example; this one is simpler (but notice the *multidimensional* columns called *pressure* and *temperature*!). And we will introduce a new manner to describe a Table as a dictionary, as you can see in the Event description. See section 4.4.2 about the different kinds of descriptor objects that can be passed to the `createTable()` method.

```
from numarray import *
from tables import *

# Describe a particle record
class Particle(IsDescription):
    name          = StringCol(length=16) # 16-character String
    lati          = IntCol()             # integer
    longi         = IntCol()             # integer
    pressure      = Float32Col(shape=(2,3)) # array of floats (single-precision)
    temperature   = FloatCol(shape=(2,3))  # array of doubles (double-precision)

# Another way to describe the columns of a table
Event = {
    "name"       : Col('CharType', 16),    # 16-character String
    "TDCcount"   : Col("UInt8", 1),        # unsigned byte
    "ADCcount"   : Col("UInt16", 1),       # Unsigned short integer
    "xcoord"     : Col("Float32", 1),      # integer
    "ycoord"     : Col("Float32", 1),      # integer
}

# Open a file in "w"rite mode
fileh = openFile("tutorial2.h5", mode = "w")
# Get the HDF5 root group
```



```
root = fileh.root
# Create the groups:
for groupname in ("Particles", "Events"):
    group = fileh.createGroup(root, groupname)
# Now, create and fill the tables in Particles group
gparticles = root.Particles
# Create 3 new tables
for tablename in ("TParticle1", "TParticle2", "TParticle3"):
    # Create a table
    table = fileh.createTable("/Particles", tablename, Particle,
                              "Particles: "+tablename)
    # Get the record object associated with the table:
    particle = table.row
    # Fill the table with 257 particles
    for i in xrange(257):
        # First, assign the values to the Particle record
        particle['name'] = 'Particle: %6d' % (i)
        particle['lati'] = i
        particle['longi'] = 10 - i
        ##### Detectable errors start here. Play with them!
        particle['pressure'] = array(i*arange(2*3), shape=(2,4)) # Incorrect
        #particle['pressure'] = array(i*arange(2*3), shape=(2,3)) # Correct
        ##### End of errors
        particle['temperature'] = (i**2) # Broadcasting
        # This injects the Record values
        particle.append()
    # Flush the table buffers
    table.flush()

# Now, go for Events:
for tablename in ("TEvent1", "TEvent2", "TEvent3"):
    # Create a table in Events group
    table = fileh.createTable(root.Events, tablename, Event,
                              "Events: "+tablename)
    # Get the record object associated with the table:
    event = table.row
    # Fill the table with 257 events
    for i in xrange(257):
        # First, assign the values to the Event record
        event['name'] = 'Event: %6d' % (i)
        event['TDCcount'] = i % (1<<8) # Correct range
        ##### Detectable errors start here. Play with them!
        #event['xcoord'] = float(i**2) # Correct spelling
        event['xcoor'] = float(i**2) # Wrong spelling
        event['ADCcount'] = i * 2 # Correct type
        #event['ADCcount'] = "s" # Wrong type
        ##### End of errors
        event['ycoord'] = float(i)**4
        # This injects the Record values
        event.append()

# Flush the buffers
table.flush()
```

```
# Read the records from table "/Events/TEvent3" and select some
table = root.Events.TEvent3
e = [ p['TDCcount'] for p in table
      if p['ADCcount'] < 20 and 4 <= p['TDCcount'] < 15 ]
print "Last record ==>", p
print "Selected values ==>", e
print "Total selected records ==> ", len(e)
# Finally, close the file (this also will flush all the remaining buffers!)
fileh.close()
```

3.3.1 Shape checking

If you have read the code carefully it looks pretty good, but it won't work. When you run this example, you will get the next error:

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 53, in ?
    particle['pressure'] = array(i*arange(2*3), shape=(2,4)) # Incorrect
  File "/usr/local/lib/python2.2/site-packages/numarray/numarraycore.py", line 281, in array
    a.setshape(shape)
  File "/usr/local/lib/python2.2/site-packages/numarray/generic.py", line 530, in setshape
    raise ValueError("New shape is not consistent with the old shape")
ValueError: New shape is not consistent with the old shape
```

which is saying that you are trying to assign an array of incompatible shape to a table cell. If you look at the source, we were trying to assign an array of shape (2,4) to a pressure element, which was defined to have a shape of (2,3).

In general, this kind of operations are forbidden, with a honorable exception: when you tries to assign an *scalar* value to a column cell that is multidimensional, all the cell elements are populated with the value of this scalar. This happens in the next line:

```
particle['temperature'] = (i**2) # Broadcasting
```

So, the value `i**2` is assigned to all the elements of the temperature table cell. This capability is provided by the numarray package and is known as *broadcasting*.

3.3.2 Field name checking

After fixing the previous error, and re-running again the program, we will get another one:

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 74, in ?
    event['xcoor'] = float(i**2) # Wrong spelling
  File "/home/falted/PyTables/pytables-0.7/src/hdf5Extension.pyx",
line 1812, in hdf5Extension.Row.__setitem__
    raise AttributeError, "Error setting \"%s\" attr.\n %s" % \
AttributeError: Error setting "xcoor" attr.
Error was: "exceptions.KeyError: xcoor"
```

This error is telling us that we tried to assign a value to a non-existent field in the *event* table object. By looking carefully at the Event class attributes, we see that we misspelled the *xcoord* field (we wrote *xcoor* instead). This is very unusual in Python because if you try to assign a value to a non-existent instance variable, a new one is created with that name. Such a feature is not satisfactory when we are dealing with an object that has fixed list of field names. So, a check is made inside PyTables so that if you try to assign a value to a non-existing field a *KeyError* is raised.

Figure 3.3 shows the PyTables GUI interface. The left pane displays the file hierarchy for 'tutorial2.h5', showing 'Events' and 'Particles' groups, each containing three tables (TEvent1, TEvent2, TEvent3 and TParticle1, TParticle2, TParticle3 respectively). The main pane displays two 'TableView' windows. The top window shows the 'TEvent2' table with columns: ADCcount, TDCcount, name, xcoord, and ycoord. The bottom window shows the 'TParticle2' table with columns: lati, longi, name, pressure, and temperature. The 'TParticle2' table contains 5 rows (0-4) and 5 columns, with the 'pressure' and 'temperature' columns containing multidimensional arrays.

Figure 3.3: Table hierarchy for tutorial 2.

3.3.3 Data type checking

Finally, in order to test the type checking, we will change the next line:

```
event.ADCcount = i * 2          # Correct type
```

to read:

```
event.ADCcount = "s"           # Wrong type
```

After this modification, the next exception will be raised when the script is executed:

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 76, in ?
    event['ADCcount'] = "s"          # Wrong type
  File "/home/faltd/PyTables/pytables-0.7/src/hdf5Extension.pyx", line 1812, in hdf5Extension.raise AttributeError, "Error setting \"%s\" attr.\n %s" % \
AttributeError: Error setting "ADCcount" attr.
Error was: "exceptions.TypeError: NA_setFromPythonScalar: bad value type."
```

that states the kind of error (TypeError).

You can admire the structure we have created with this (corrected) script in figure 3.3. In particular, pay attention to the multidimensional column cells in table /Particles/TParticle2.

Feel free to visit the rest of examples in directory examples, and try to understand them. I've tried to make several use cases to give you an idea of the PyTables capabilities and its way of dealing with HDF5 objects.

*"Tenho pensamentos que, se pudesse
revelá-los e fazê-los viver, acrescentariam
nova luminosidade às estrelas, nova
beleza ao mundo e maior amor ao
coração dos homens."*

—Fernando Pessoa, in "O Eu Profundo"

Chapter 4

Library Reference

PyTables implements several classes to represent the different nodes in the object tree. They are named `File`, `Group`, `Leaf`, `Table` and `Array`. Another one is responsible to build record objects from a subclass user declaration, and performs field and type checks; its name is `IsDescription`. An important function, called `openFile` is responsible to create, open or append to files. In addition, a few utility functions are defined to guess if the user supplied file is a PyTables or HDF5 file. These are called `isPyTablesFile` and `isHDF5`. Finally, several first-level variables are also available to the user that informs about PyTables version, file format version or underlying libraries (as for example HDF5) version number.

Let's start discussing the first-level variables and functions available to the user, then the methods in the classes defined in PyTables.

4.1 tables variables and functions

4.1.1 Global variables

__version__ The PyTables version number.

HDF5Version The underlying HDF5 library version number.

ExtVersion The Pyrex extension types version. This might be useful when reporting bugs.

4.1.2 Global functions

openFile(filename, mode='r', title='', trMap={}, rootUEP='/') Open a PyTables file and returns a File object.

filename The name of the file (supports environment variable expansion). It is suggested that it should have any of ".h5", ".hdf" or ".hdf5" extensions, although this is not mandatory.

mode The mode to open the file. It can be one of the following:

'r' read-only; no data can be modified.

'w' write; a new file is created (an existing file with the same name is deleted).

'a' append; an existing file is opened for reading and writing, and if the file does not exist it is created.

'r+' is similar to 'a', but the file must already exist.

title If filename is new, this will set a title for the root group in this file. If filename is not new, the title will be read from disk, and this will not have any effect.

trMap A dictionary to map names in the object tree Python namespace into different HDF5 names in file namespace. The keys are the Python names, while the values are the HDF5 names. This is useful when you need to name HDF5 nodes with invalid or reserved words in Python.

rootUEP The root User Entry Point. This is a group in the HDF5 hierarchy which will be taken as the starting point to create the object tree. The group has to be named after its HDF5 name and can be a path. If it does not exist, a `RuntimeError` is issued. Use this if you do not want to build the **entire** object tree, but rather only a **subtree**.

isHDF5(filename) Determines whether filename is in the HDF5 format or not. When successful, returns a positive value, for TRUE, or 0 (zero), for FALSE. Otherwise returns a negative value. To this function to work, it needs a closed file.

isPyTablesFile(filename) Determines whether a file is in the PyTables format. When successful, returns the format version string, for TRUE, or 0 (zero), for FALSE. Otherwise returns a negative value. To this function to work, it needs a closed file.

4.2 The `IsDescription` class

This class is in fact a so-called *metaclass* object. There is nothing special on this fact, except that their subclasses attributes are transformed during its instantiation phase, and new methods for instances are defined based on the values of the class attributes.

It is designed to be used as an easy, yet meaningful way to describe the properties of `Table` objects through the use of classes that inherit properties from it. In order to define such a special class, you have to declare it as descendant of `IsDescription`, with many attributes as columns you want in your table. The name of these attributes will become the name of the columns, while its values are the properties of the columns that are obtained through the use of the `Col` class constructor. See the section 4.3 for instructions on how to define the properties of the table columns.

Then, you can pass an instance of this object to the `Table` constructor, where all the information it contains will be used to define the table structure. See the section 3.3 for an example on how that works.

4.3 The `Col` class and its descendants

The `Col` class is used as a mean to declare the different properties of a column of a table. In addition, a series of descendant classes are offered in order to make these column descriptions easier to the user. In general, it is recommended to use these descendants classes, as they are meaningful when found in the middle of the code.

The only public method accessible in these classes is the constructor itself.

`Col(dtype="Float64", shape=1, dflt=None, pos=None)` Define properties for a `Table` column.

`dtype` The data type for the column. See the appendix A for a relation of data types supported in a `IsDescription` class declaration. The type description is accepted both in string format and as `numpy` data type.

`shape` An integer or a tuple, that specifies the number of `dtype` items for each element (or shape, for multidimensional elements) of this column. For `CharType` columns, the first dimension is used as the length of the character strings. For this kind of objects, the use of `StringCol` subclass is recommended.

`dflt` The default value for elements of this column. If the user does not supply a value for an element while filling a table, this default value will be written to disk. If the user supplies a scalar value for a multidimensional column, this value is automatically *broadcasted* to all the elements in the column cell. If `dflt` is not supplied, an appropriate zero value (or *null* string) will be chosen by default.

`pos` By default, columns are disposed in memory following an alphanumerical order of the column names. In some situations, however, it is convenient to impose a user defined ordering. `pos` parameter allows the user to force the wanted disposition.

StringCol(length=None, dflt=None, shape=1, pos=None) Define a column to be of `CharType` type. The `length` parameter sets the length of the strings. The meaning of the other parameters are like in the `Col` class.

IntCol(dflt=0, shape=1, itemsize=4, sign=1, pos=None) Define a column to be of `IntXXType` type, depending on the value of `itemsize`. The `itemsize` parameter sets the number of bytes of the integers in the column and the default is 4 bytes. `sign` determines if the integers are signed or not. The meaning of the other parameters are like in the `Col` class.

This class has several descendants:

Int8Col(dflt=0, shape=1, pos=None) Define a column as an `Int8` type.

UInt8Col(dflt=0, shape=1, pos=None) Define a column as an `UInt8` type.

Int16Col(dflt=0, shape=1, pos=None) Define a column as an `Int16` type.

UInt16Col(dflt=0, shape=1, pos=None) Define a column as an `UInt16` type.

Int32Col(dflt=0, shape=1, pos=None) Define a column as an `Int32` type.

UInt32Col(dflt=0, shape=1, pos=None) Define a column as an `UInt32` type.

Int64Col(dflt=0, shape=1, pos=None) Define a column as an `Int64` type.

UInt64Col(dflt=0, shape=1, pos=None) Define a column as an `UInt64` type.

FloatCol(dflt=0, shape=1, itemsize=8, pos=None) Define a column to be of `FloatXXType` type, depending on the value of `itemsize`. The `itemsize` parameter sets the number of bytes of the floats in the column and the default is 8 bytes (double precision). The meaning of the other parameters are like in the `Col` class.

This class has two descendants:

Float32Col(dflt=0.0, shape=1, pos=None) Define a column as a `Float32` type.

Float64Col(dflt=0.0, shape=1, pos=None) Define a column as a `Float64` type.

4.4 The File class

This class is returned when a `PyTables` file is opened with the `openFile()` function. It has methods to create, open, flush and close `PyTables` files. Also, `File` class offer methods to traverse the object tree, as well as to create, rename and delete nodes. One of its attributes (`rootUEP`) represents the *user entry point* to the object tree attached to the file.

Next, we will discuss the attributes and methods for `File` class¹.

4.4.1 File instance variables

filename Filename opened.

isopen It takes the value 1 if the underlying file is open. 0 otherwise.

mode Mode in which the filename was opened.

title The title of the root group in file.

rootUEP The UEP (User Entry Point) group in file (see 4.1.2).

trMap This is a dictionary that maps node names between python and HDF5 domain names. Its initial values are set from the `trMap` parameter passed to the `openFile()` function. You can change its contents after a file is opened and the new map will take effect over any new object added to the tree.

objects Dictionary with all objects (groups or leaves) on tree.

groups Dictionary with all object groups on tree.

leaves Dictionary with all object leaves on tree.

¹ On the following, the term `Leaf` will whether refer to a `Table` or `Array` node object.

4.4.2 File methods

createGroup(*where*, *name*, *title*=*"*)

Create a new Group instance with name *name* in *where* location.

where The parent group where the new group will hang. *where* parameter can be a path string (for example `"/Particles/TParticle1"`), or another Group instance.

name The name of the new group.

title A description for this group.

createTable(*where*, *name*, *description*, *title*=*"*, *compress*=0, *complib* = *'zlib'*, *expectedrows*=10000)

Create a new Table instance with name *name* in *where* location.

where The parent group where the new table will hang. *where* parameter can be a path string (for example `"/Particles/TParticle1"`), or Group instance.

name The name of the new table.

description An instance of a user-defined class (derived from the `IsDescription` class) where table fields are defined. However, in certain situations, it is more handy to allow this description to be supplied as a dictionary (for example, when you do not know beforehand which structure will have your table). In such a cases, you can pass the description as a dictionary as well. See section 3.3 for an example of use. Finally, a `RecArray` object from the `numarray` package is also accepted, and all the information about columns and other metadata is used as a basis to create the Table object. Moreover, if the `RecArray` has actual data this is also injected on the newly created Table object.

title A description for this object.

compress Specifies a compress level for data. The allowed range is 0-9. A value of 0 disables compression. The default is that compression is disabled, that balances between compression effort and CPU consumption.

complib Specifies the compression library to be used. Right now, `"zlib"` (default), `"lzo"` and `"ucl"` values are supported. See section 5.2 for some advice on which library is better suited to your needs.

expectedrows An user estimate of the number of records that will be on table. If not provided, the default value is appropriate for tables until 1 MB in size (more or less, depending on the record size). If you plan to save bigger tables you should provide a guess; this will optimize the HDF5 B-Tree creation and management process time and memory used. See section 5.3 for a detailed justification of that issue.

createArray(*where*, *name*, *object*, *title*=*"*)

Create a new Array instance with name *name* in *where* location.

where The parent group where the new array will hang. *where* parameter can be a path string (for example `"/Particles/TParticle1"`), or Group instance.

name The name of the new array.

object The regular array to be saved. Currently accepted values are: lists, tuples, scalars (int and float), strings and (multidimensional) `Numeric` and `NumArray` arrays (including `CharArrays` string arrays). However, these objects must be regular (i.e. they cannot be like, for example, `[[1,2],2]`). Also, objects that has some of its dimension equal to zero are not supported (this will be solved when unlimited arrays will be implemented).

title A description for this object.

getNode(where, name="", classname="")

Returns the object node *name* under *where* location.

where Can be a path string or Group instance. If *where* doesn't exist or has not a child called *name*, a `ValueError` error is raised.

name The object name desired. If *name* is a null string (''), or not supplied, this method assumes to find the object in *where*.

classname If supplied, returns only an instance of this class name. Allowed names in *classname* are: 'Group', 'Leaf', 'Table' and 'Array'. Note that these values are strings.

getAttrNode(where, attrname, name="")

Returns the attribute *attrname* under *where.name* location.

where Can be a path string or Group instance. If *where* doesn't exist or has not a child called *name*, a `ValueError` error is raised.

attrname The name of the attribute to get.

name The node name desired. If *name* is a null string (''), or not supplied, this method assumes to find the object in *where*.

setAttrNode(where, attrname, attrvalue, name="")

Sets the attribute *attrname* with value *attrvalue* under *where.name* location.

where Can be a path string or Group instance. If *where* doesn't exist or has not a child called *name*, a `ValueError` error is raised.

attrname The name of the attribute to set on disk.

attrvalue The value of the attribute to set. Only strings attributes are supported natively right now. However, you can always use `(c)Pickle` so as to serialize any object you want save therein.

name The node name desired. If *name* is a null string (''), or not supplied, this method assumes to find the object in *where*.

listNodes(where, classname="")

Returns a list with all the object nodes (Group or Leaf) hanging from *where*. The list is alphanumerically sorted by node name.

where The parent group. Can be a path string or Group instance.

classname If a *classname* parameter is supplied, the iterator will return only instances of this class (or subclasses of it). The only supported classes in *classname* are 'Group', 'Leaf', 'Table' and 'Array'. Note that these values are strings.

removeNode(where, name = "", recursive=0)

Removes the object node *name* under *where* location.

where Can be a path string or Group instance. If *where* doesn't exist or has not a child called *name*, a `LookupError` error is raised.

name The name of the node to be removed. If not provided, the *where* node is changed.

recursive If not supplied, the object will be removed only if it has no children. If supplied with a true value, the object and all its descendants will be completely removed.

renameNode(where, newname, name)

Rename the object node *name* under *where* location.

where Can be a path string or Group instance. If *where* doesn't exist or has not a child called *name*, a `LookupError` error is raised.

newname Is the new name to be assigned to the node.

name The name of the node to be changed. If not provided, the *where* node is changed.

walkGroups(where='/')

Iterator that returns the list of Groups (not Leaves) hanging from *where*. If *where* is not supplied, the root object is taken as origin. The returned Group list is in a top-bottom order, and alphanumerically sorted when they are at the same level.

where The origin group. Can be a path string or Group instance.

flush()

Flush all the leaves in the object tree.

close()

Flush all the leaves in object tree and close the file.

4.4.3 File special methods

Following are described the methods that automatically trigger actions when a `File` instance is accessed in a special way (e.g., `fileh("/detector")` will cause a call to `group.__call__("/detector")`).

__call__(where="/", classname="")

Recursively iterate over the children in the `File` instance. It takes two parameters:

where If supplied, the iteration starts from this group.

classname (*String*) If supplied, only instances of this class are returned.

Example of use:

```
# Recursively print all the nodes hanging from '/detector'
print "Nodes hanging from group '/detector':"
for node in h5file("/detector"):
    print node
```

__iter__()

Iterate over the children on the `File` instance. However, this does not accept parameters. This iterator is *recursive*.

Example of use:

```
# Recursively list all the nodes in the object tree
print "All nodes in the object tree:"
for node in h5file:
    print node
```

4.5 The Group class

Instances of this class are a grouping structure containing instances of zero or more groups or leaves, together with supporting metadata.

Working with groups and leaves is similar in many ways to working with directories and files, respectively, in a Unix filesystem. As with Unix directories and files, objects in the object tree are often described by giving their full (or absolute) path names. This full path can be specified either as string (like in `'/group1/group2'`) or as a complete object path written in the Pythonic fashion known as *natural name* schema (like in `file.root.group1.group2`) and discussed in the section 1.2.

A collateral effect of the *natural naming* schema is that you must be aware when assigning a new attribute variable to a Group object to not collide with existing children node names. For this reason and to not pollute the children namespace, it is explicitly forbidden to assign "normal" attributes to Group instances, and the only ones allowed must start with some reserved prefixes, like `"_f_"` (for methods) or `"_v_"` (for instance variables) prefixes. Any attempt to assign a new attribute that does not starts with these prefixes, will raise a `NameError` exception.

Other effect is that you cannot use reserved Python names or other non-allowed python names (like for example `"$a"` or `"44"`) as node names. You can, however, make use of a translation map dictionary in the `File.openfile()` method (see section 4.1.2) so as to use non valid Python names as node names in the file.

4.5.1 Group instance variables

`_v_title` A description for this group.

`_v_name` The name of this group.

`_v_hdf5name` The name of this group in HDF5 file namespace.

`_v_pathname` A string representation of the group location in tree.

`_v_parent` The parent Group instance.

`_v_rootgroup` Pointer to the root group object.

`_v_file` Pointer to the associated File object.

`_v_childs` Dictionary with all nodes (groups or leaves) hanging from this instance.

`_v_groups` Dictionary with all node groups hanging from this instance.

`_v_leaves` Dictionary with all node leaves hanging from this instance.

`_v_attrs` The associated `AttributeSet` instance (see 4.10).

4.5.2 Group methods

This class define the `__setattr__`, `__getattr__` and `__delattr__` and they work as normally intended. So, you can access, assign or delete childs to a group by just using the next constructs:

```
# Add a Table child instance under group with name "tablename"
group.tablename = Table(recordDict, "Record instance")
table = group.tablename      # Get the table child instance
del group.tablename          # Delete the table child instance
```

Caveat: The following methods are documented for completeness, and they can be used without any problem. However, you should use the high-level counterpart methods in the `File` class, because these are most used in documentation and examples, and are a bit more powerful than ones those exposed here.

`_f_join(name)` Helper method to correctly concatenate a name child object with the pathname of this group.

`_f_rename(newname)` Change the name of this group to *newname*.

`_f_remove(recursive=0)` Remove this object. If *recursive* is true, force the removal even if this group has children.

`_f_getAttr(attrname)` Gets the HDF5 attribute *attrname* of this group.

`_f_setAttr(attrname, attrvalue)` Sets the attribute *attrname* of this group to the value *attrvalue*. Only string values are allowed.

`_f_listNodes(classname='')` Returns a *list* with all the object nodes hanging from this instance. The list is alphanumerically sorted by node name. If a *classname* parameter is supplied, it will only return instances of this class (or subclasses of it). The supported classes in *classname* are 'Group', 'Leaf', 'Table' and 'Array'.

`_f_walkGroups()` Iterator that returns the list of Groups (not Leaves) hanging from *self*. The returned Group list is in a top-bottom order, and alphanumerically sorted when they are at the same level.

`_f_close()` Close this group, making it and its children unaccessible in the object tree.

4.5.3 Group special methods

Following are described the methods that automatically trigger actions when a Group instance is accessed in a special way (e.g., `group("Table")` will cause a call to `group.__call__("Table")`).

`__call__(classname="", recursive=0)`

Iterate over the childs in the Group instance. It takes two parameters:

classname (*String*) If supplied, only instances of this class are returned.

recursive (*Integer*) If false, only childs hanging immediately after the group are returned. If true, a recursion over all the groups hanging from it is performed.

Example of use:

```
# Recursively print all the arrays hanging from '/'
print "Arrays the object tree '':"
for array in h5file.root(classname="Array", recursive=1):
    print array
```

`__iter__()`

Iterate over the childs on the group instance. However, this does not accept parameters. This iterator is not recursive.

Example of use:

```
# Non-recursively list all the nodes hanging from '/detector'
print "Nodes in '/detector' group:"
for node in h5file.root.detector:
    print node
```

4.6 The Leaf class

This is a helper class useful to place common functionality of all Leaf objects. It is also useful for classifying purposes. A Leaf object is an end-node, that is, a node that can hang directly from a group object, but that is not a group itself. Right now this set is composed by `Table` and `Array` objects. In fact, `Table` and `Array` classes inherit functionality from this class using the *mix-in* technique.

The public variables and methods that `Table` and `Array` inherits from `Leaf` are listed below.

4.6.1 Leaf instance variables

name The Leaf node name in Python namespace.

hdf5name The Leaf node name in HDF5 namespace.

title The Leaf title.

shape The shape of the associated data in the Leaf.

byteorder The byteorder of the associated data of the Leaf.

attrs The associated `AttributeSet` instance (see 4.10).

4.6.2 Leaf methods

rename(newname) Change the name of this leaf to *newname*.

remove() Remove this leaf.

getAttr(attrname) Gets the HDF5 attribute *attrname* of this leaf.

setAttr(attrname, attrvalue) Sets the attribute *attrname* of this leaf to the value *attrvalue*. Only string values are allowed.

flush() Flush the leaf buffers.

close() Flush the leaf buffers and close the HDF5 dataset.

4.7 The Table class

Instances of this class represents table objects in the object tree. It provides methods to create new tables or open existing ones, as well as methods to read/write data and metadata from/to table objects in the file.

Data can be read from or written to tables by accessing to an special object that hangs from `Table`. This object is an instance of the `Row` class (see 4.8). See the tutorial sections chapter 3 on how to use the `Row` interface.

Please note that this object inherits all the public attributes and methods that `Leaf` has.

4.7.1 Table instance variables

description The metaobject describing this table

row The `Row` instance for this table (see 4.8).

nrows The number of rows in this table.

colnames The field names for the table (list).

coltypes The data types for the table fields (dictionary).

colshapes The shapes for the table fields (dictionary).

4.7.2 Table methods

iterrows(start=None, stop=None, step=None)

Returns an iterator yielding Row instances built from rows in table. If a range is supplied (i.e. some of the *start*, *stop* or *step* parameters are passed), only the appropriate rows are returned. Else, all the rows are returned.

start Sets the starting row to return data. A value of 0 means the first row. It accepts negative values meaning that the count starts from the end.

stop Sets the last row to be returned to *stop* - 1, i.e. the end point is omitted (in the Python range tradition). It accepts, likewise *start*, negative values. A special value of 0 means the last row.

step When step is given, it specifies the increment. Negative values are not allowed right now.

read(self, start=None, stop=None, step=None, field=None, flavor=None)

Returns the actual data in Table. If *field* is not supplied, it returns the data as a RecArray object table.

start Sets the starting row to return data. A value of 0 means the first row. It accepts negative values meaning that the count starts from the end.

stop Sets the last row to be returned to *stop* - 1, i.e. the end point is omitted (in the Python range tradition). It accepts, likewise *start*, negative values. A special value of 0 means the last row.

step When step is given, it specifies the increment. Negative values are not allowed right now.

field If specified, only the column *field* is returned as a NumArray object. If this is not supplied, all the fields are selected and a RecArray is returned.

flavor When a field in table is selected, passing a *flavor* parameter make an additional conversion to happen in the default NumArray object. *flavor* must have any of the next values: Numeric, Tuple or List.

removeRows(start=None, stop=None)

Removes a range of rows in the table. If only *start* is supplied, this row is to be deleted. If a range is supplied, i.e. both the *start* and *stop* parameters are passed, all the rows in the range are removed². A *step* parameter is not supported yet.

start Sets the starting row to be removed. It accepts negative values meaning that the count starts from the end. A value of 0 means the first row.

stop Sets the last row to be removed to *stop* - 1, i.e. the end point is omitted (in the Python range tradition). It accepts, likewise *start*, negative values. A special value of 0 means the last row.

4.7.3 Table special methods

Following are described the methods that automatically trigger actions when a Table instance is accessed in a special way (e.g., `table["var2"]` will cause a call to `table.__getitem__("var2")`).

² However, for `removeRows()` to work, you need that the rows **after** the *stop* parameter will fit in-memory so as to method to work. This limitation will be hopefully removed in a future version.

__call__(start=None, stop=None, step=None)

It returns the same iterator than `Table.iterrows(start, stop, step)`. It is, therefore, a shorter way to call it.

Example of use:

```
result = [ row['var2'] for row in table(step=4)
           if row['var1'] <= 20 ]
```

Which is equivalent to:

```
result = [ row['var2'] for row in table.iterrows(step=4)
           if row['var1'] <= 20 ]
```

__iter__()

It returns the same iterator than `Table.iterrows(0,0,1)`. However, this does not accept parameters.

Example of use:

```
result = [ row['var2'] for row in table
           if row['var1'] <= 20 ]
```

Which is equivalent to:

```
result = [ row['var2'] for row in table.iterrows()
           if row['var1'] <= 20 ]
```

__getitem__(key)

It takes different actions depending on the type of the key parameter:

key is an Integer The corresponding table row is returned as a `RecArray.Record` object.

key is a Slice The row slice determined by key is returned as a `RecArray` object.

key is a String The key is interpreted as a *column* name of the table, and, if it exists, it is read and returned as a `NumArray` or `CharArray` object (whatever is appropriate).

Example of use:

```
record = table[4]
recarray = table[4:1000:2]
narray = table["var2"]
```

Which is equivalent to:

```
record = table.read(start=4)[0]
recarray = table.read(start=4, stop=1000, step=2)
narray = table.read(field="var2")
```

4.8 The Row class

This class is used to fetch and set values on the table fields. It works very much like a dictionary, where the keys are the field names of the associated table and the values are the values of those fields in a specific row.

This object turns out to actually be an extension type, so you won't be able to access their documentation interactively. Neither you won't be able to access its internal attributes (they are not directly accessible from Python), although that *accessors* (i.e. methods that return an internal attribute) has been defined for the most important variables.

4.8.1 Row methods

append() Once you have filled the proper fields for the current row, calling this method actually commit this data to the disk (actually data is written to the output buffer).

nrow() Accessor that returns the current row in the table. It is useful to know which row is being dealt with in the middle of a loop.

4.9 The Array class

Represents an array on file. It provides methods to create new arrays or open existing ones, as well as methods to write/read data and metadata to/from array objects in the file.

Caveat: All `Numeric` and `numarray` data types are supported except those that corresponds to complex data types³. See `numarray` manual (Greenfield *et al.*) to know more about the supported data types, or see appendix A.

Please note that this object inherits all the public attributes and methods from `Leaf`.

4.9.1 Array instance variables

type The type class of the represented array.

flavor The string object representation for this array. It can be any of `"NumArray"`, `"CharArray"`, `"Numeric"`, `"List"`, `"Tuple"`, `"String"`, `"Int"` or `"Float"` values.

4.9.2 Array methods

Note that, as this object has not internal I/O buffers, there is no point in calling `flush()` method inherited from `Leaf`.

read() Read the array from disk and return it as a `NumArray` (default) object, or if possible, with the original *flavor* that it was saved. The supported flavors are: `NumArray`, `CharArray`, `Numeric`, `List`, `Tuple`, `String`, `Int` or `Float`. Note that as long as this method is not called, the actual array data is resident on disk, not in memory.

4.10 The AttributeSet class

Represents the set of attributes of a node (`Leaf` or `Group`). It provides methods to create new attributes, open, rename or delete existing ones.

Like in `Group` instances, `AttributeSet` instances use a special feature called *natural naming*, i.e. you can access the attributes on disk like if they were *normal* `AttributeSet` attributes. This offers the user a very convenient way to access (but also set and delete) node attributes by simply specifying them like a *normal* attribute class.

Caveat: All Python data types are supported. The scalar ones (i.e. `String`, `Int` and `Float`) are mapped directly to the HDF5 counterparts, so you can correctly visualize them with any HDF5 tool. However, the

³ However, these might be included in the future

rest of the data types and more general objects are serialized using `cPickle`, so you will be able to correctly retrieve them only from a Python-aware HDF5 library. Hopefully, the list of supported native attributes will be extended to multidimensional arrays sometime in the future.

4.10.1 AttributeSet instance variables

`_v_node` The parent node instance.

`_v_attrnames` List with all attribute names.

`_v_attrnamesys` List with system attribute names.

`_v_attrnamesuser` List with user attribute names.

4.10.2 AttributeSet methods

Note that this class define the `__setattr__`, `__getattr__` and `__delattr__` and they work as normally intended. So, you can access, assign or delete attributes on disk by just using the next constructs:

```
leaf.attrs.myattr = "string attr" # Set the attribute myattr
attrib = leaf.attrs.myattr # Get the attribute myattr
del leaf.attrs.myattr # Delete the attribute myattr
```

`_f_list(attrset = "user")` Return the list of attributes of the parent node.

`attrset` Selects the attribute set to be returned. An "user" value returns only the user attributes. This is the default. "sys" returns only the system (some of which are read-only) attributes. "readonly" returns the system read-only attributes. "all" returns both the system and user attributes.

`_f_rename(oldattrname, newattrname)` Rename an attribute.

... durch planmässiges Tattonieren.

*—Johann Karl Friedrich Gauss
[asked how he came upon his theorems]*

Chapter 5

Optimization tips

On this chapter, you will get deeper knowledge of PyTables internals. PyTables has several places where the user can improve the performance of his application. If you are planning to deal with really large data, you should read carefully this section in order to learn how to get an important boost for your code. But if your dataset is small or medium size (say, up to 1 MB), you should not worry about that as the default parameters in PyTables are already tuned to handle that perfectly.

5.1 Taking advantage of Psyco

Psyco (see Rigo) is a kind of specialized compiler for Python that typically accelerates Python applications with no change in source code. You can think of Psyco as a kind of just-in-time (JIT) compiler, a little bit like Java's, that emit machine code on the fly instead of interpreting your Python program step by step. The result is that your unmodified Python programs run faster.

Psyco is very easy to install and use, so in most scenarios it is worth to have it a try. However, it only runs on Intel 386 architectures, so if you are using other architectures, you are out of luck (at least until Psyco will support yours).

As an example, imagine that you have a small script that reads and selects data over a series of datasets, like this:

```
def readFile(filename):
    "Select data from all the tables in filename"

    fileh = openFile(filename, mode = "r")
    result = []
    for table in fileh("/", 'Table'):
        result = [ p['var3'] for p in table if p['var2'] <= 20 ]

    fileh.close()
    return e

if __name__=="__main__":
    print readFile("myfile.h5")
```

In order to accelerate this piece of code, you can rewrite your main program to look like:

```
if __name__=="__main__":
    import pysco
    pysco.bind(readFile)
    print readFile("myfile.h5")
```

That's all!. From now on, each time that you execute your python script, Psyco will deploy its sophisticated algorithms so as to accelerate your calculations.

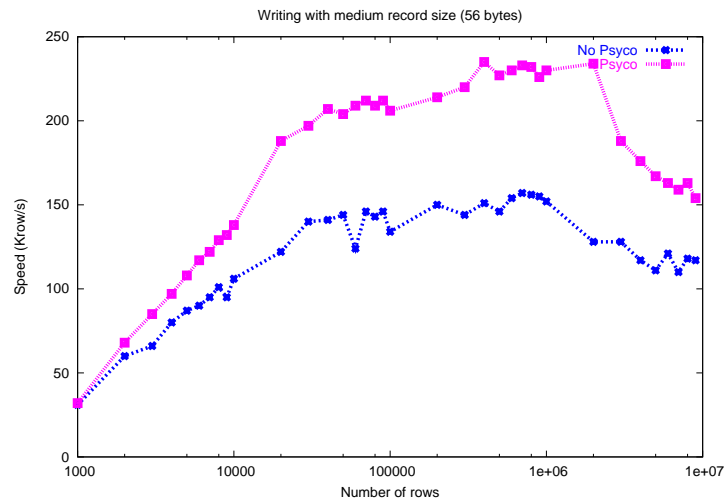


Figure 5.1: Writing tables with/without Psyco.

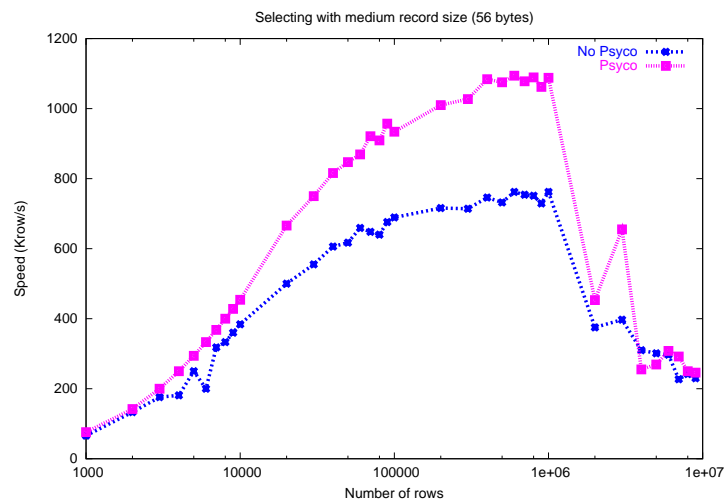


Figure 5.2: Reading tables with/without Psyco.

You can see in the graphs 5.1 and 5.2 how much I/O speed improvement you can get by using Psyco. By looking at this figures you can get an idea if these improvements are of your interest or not. In general, if you are not going to use compression you will take advantage of Psyco if your tables are medium sized ($1e+3 < \text{nrows} < 1e+6$), and this advantage will disappear progressively when the number of rows grows well over one million. However if you use compression, you will probably see improvements even beyond this limit (see section 5.2). As always, there is no substitute for experimentation with your own dataset.

5.2 Compression issues

One of the beauties of `PyTables` is that it supports compression on tables (but not on arrays!, that may come later), although it is disabled by default. Compression of big amounts of data might be a bit controversial feature, because compression has a legend of being a very big CPU time resources consumer. However, if you are willing to check if compression can help not only reducing your dataset file size but **also** improving your I/O efficiency, keep reading.

There is an usual scenario where users need to save duplicated data in some record fields, while the others

Table 5.1: Comparison between different compression libraries. The tests has been conducted on a Pentium 4 at 2 GHz and a hard disk at 4200 RPM.

Compr. Lib	File size (MB)	Time writing (s)	Time reading (s)	Speed writing (Krow/s)	Speed reading (Krow/s)
NO COMPR	244.0	24.4	16.0	18.0	27.8
Zlib (lvl 1)	8.5	17.0	3.11	26.5	144.4
Zlib (lvl 6)	7.1	20.1	3.10	22.4	144.9
Zlib (lvl 9)	7.2	42.5	3.10	10.6	145.1
LZO (lvl 1)	9.7	14.6	1.95	30.6	230.5
UCL (lvl 1)	6.9	38.3	2.58	11.7	185.4

have varying values. In a relational database approach such a redundant data can normally be moved to other tables and a relationship between the rows on the separate tables can be created. But that takes analysis and implementation time, and made the underlying libraries more complex and slower.

PyTables transparent compression allows the user to not worry about finding which is their optimum data tables strategy, but rather use less, not directly related, tables with a larger number of columns while still not cluttering the database too much with duplicated data (compression is responsible to avoid that). As a side effect, data selections can be made more easily because you have more fields available in a single table, and they can be referred in the same loop. This process may normally end in a simpler, yet powerful manner to process your data (although you should still be careful about what kind of scenarios compression use is convenient or not).

The compression library used by default is the **Zlib** (see Gailly and Adler), and as HDF5 *requires* it, you can safely use it and expect that your HDF5 files can be read on any other platform that has HDF5 libraries installed. Zlib provides good compression ratio, although somewhat slow, and reasonably fast decompression. Because of that, it is a good candidate to be used for compress you data.

However, in many situations (i.e. write *once*, read *multiple*), it is critical to have *very good* decompression speed (at expense of whether less compression or more CPU wasted on compression, as we will see soon). This is why support for two additional compressors has been added to PyTables: LZO and UCL (see Oberhumer). Following his author (and checked by the author of this manual), LZO offers pretty fast compression (although small compression ratio) and extremely fast decompression while UCL achieve an excellent compression ratio (at the price of spending much more CPU time) while allowing very fast decompression (and *very close* to the LZO one). In fact, LZO and UCL are so fast when decompressing that, in general (that depends on your data, of course), writing and reading a compressed table is actually faster (and sometimes **much faster**) than if it is uncompressed. This fact is very important, specially if you have to deal with very large amounts of data.

Be aware that the LZO and UCL support in PyTables is not standard on HDF5, so if you are going to use your PyTables files in other contexts different from PyTables you will not be able to read them.

In order to give you a raw idea of what ratios would be achieved, and what resources would be consumed, look at the table 5.1. This table has been obtained from synthetic data and with a somewhat outdated PyTables version (0.5), so take this just as a guide because your mileage will probably vary. Have also a look at the graphs 5.3 and 5.4 (these graphs has been obtained with tables with different row sizes and PyTables version than the previous example, so, do not try to directly compare the figures). They show how evolves the speed of writing/reading rows as the size (the row number) of tables grows. Even though in these graphs the size of one single row is 56 bytes, you can most probably extrapolate this figures to other row sizes. If you are curious how well can perform compression together with Psycho, look at the graphs 5.5 and 5.6. As you can see, the results are pretty interesting.

By looking at graphs, you can expect that, generally speaking, LZO would be the fastest both compressing and uncompressing, but the one that achieves the worse compression ratio (although that may be just ok for many situations). UCL is the slowest when compressing, but is faster than Zlib when decompressing, and, besides, it achieves very good compression ratios (generally better than Zlib). Zlib represents a balance between them: it's somewhat slow compressing, the slowest during decompressing, but it normally achieves fairly good compression ratios.

So, if your ultimate goal is reading as fast as possible, choose LZO. If you want to reduce as much as

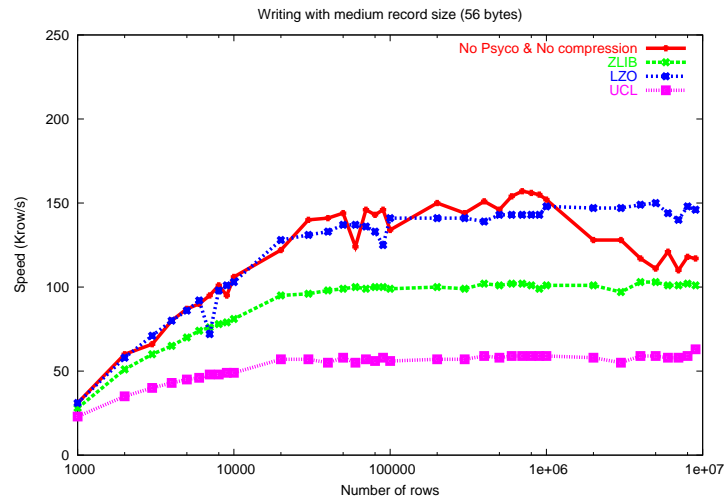


Figure 5.3: Writing tables with several compressors.

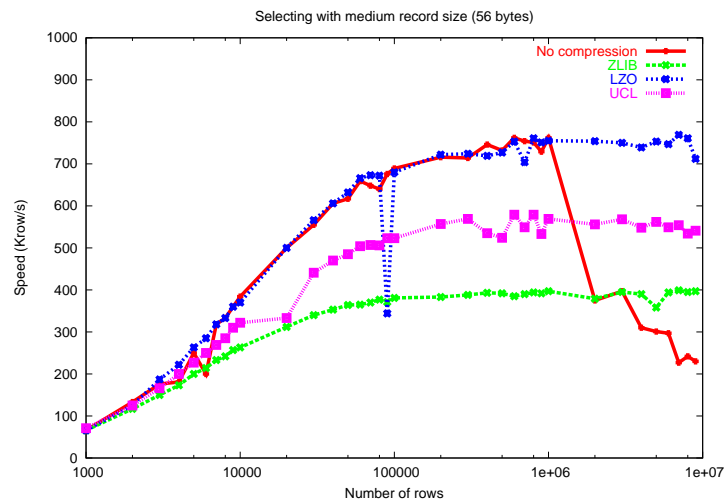


Figure 5.4: Reading tables with several compressors.

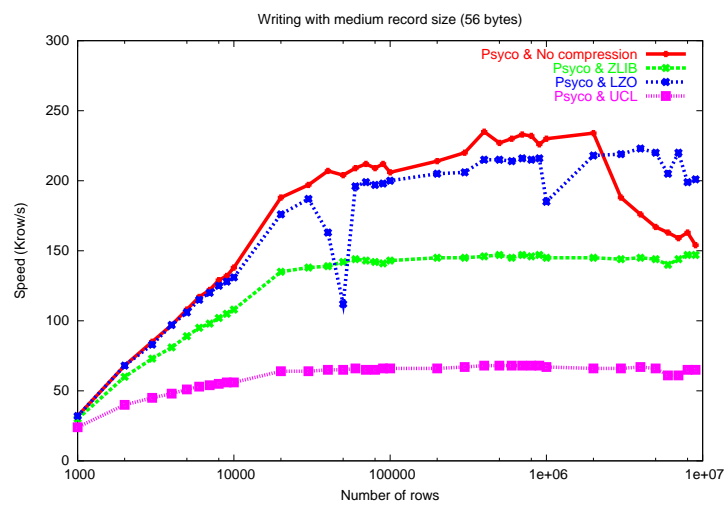


Figure 5.5: Writing tables with several compressors and Psycos.

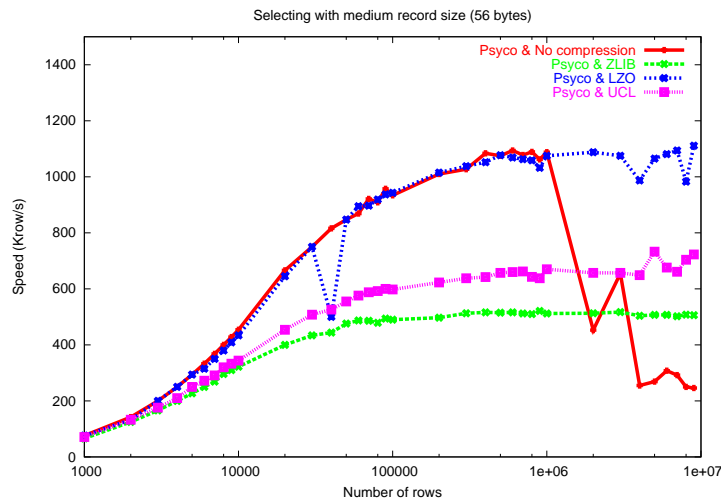


Figure 5.6: Reading tables with several compressors and Psyco.

possible your data, while retaining good read speed, choose UCL. If you don't mind too much about the above parameters and/or portability is important for you, Zlib is your best bet.

The compression level that I recommend to use for all compression libraries is 1. This is the lowest level of compression, but if you take the approach suggested above, normally the redundant data is to be found in the same row, so the redundant data locality is very high and such a small level of compression should be enough to achieve a good compression ratio on your data tables, saving CPU cycles for doing other things. Nonetheless, in some situations you may want to check how compression level affects your application.

You can select the compression library and level by setting the `complib` and `compress` keywords in the `createTable` method (see 4.4.2). A compression level of 0 will completely disable compression (the default), 1 is the less CPU time demanding level, while 9 is the maximum level and most CPU intensive. Finally, have in mind that LZO is not accepting a compression level right now, so, when using LZO, 0 means that compression is not active, and any other value means that LZO is active.

5.3 Informing PyTables about expected number of rows in tables

The underlying HDF5 library that is used by PyTables takes the data in bunches of a certain length, so-called *chunks*, to write them on disk as a whole, i.e. the HDF5 library treats chunks as atomic objects and disk I/O is always made in terms of complete chunks. This allows data filters to be defined by the application to perform tasks such as compression, encryption, checksumming, etc. on entire chunks.

An in-memory B-tree is used to map chunk structures on disk. The more chunks that are allocated for a dataset the larger the B-tree. Large B-trees take memory and causes file storage overhead as well as more disk I/O and higher contention for the metadata cache. Consequently, it's important to balance between memory and I/O overhead (small B-trees) and time to access to data (big B-trees).

PyTables can determine an optimum chunk size to make B-trees adequate to your dataset size if you help it by providing an estimation of the number of rows for a table. This must be made in table creation time by passing this value in the `expectedrows` keyword of `createTable` method (see 4.4.2).

When your table size is bigger than 1 MB (take this figure only as a reference, not strictly), by providing this guess of the number of rows you will be optimizing the access to your data. When the table size is larger than, say 100MB, you are **strongly** suggested to provide such a guess; failing to do that may cause your application doing very slow I/O operations and demanding **huge** amounts of memory. You have been warned!

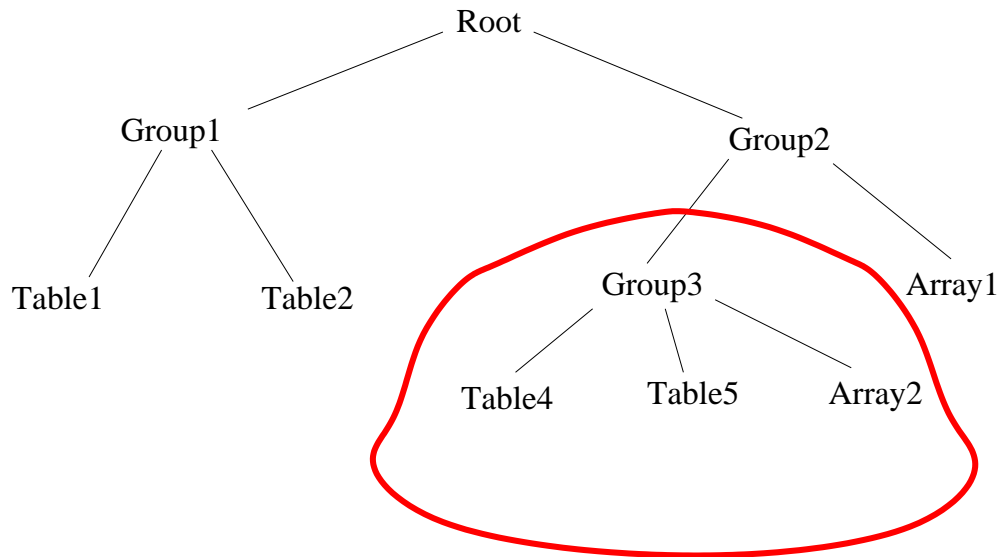


Figure 5.7: Complete tree in file `test.h5`, and subtree of interest for the user.

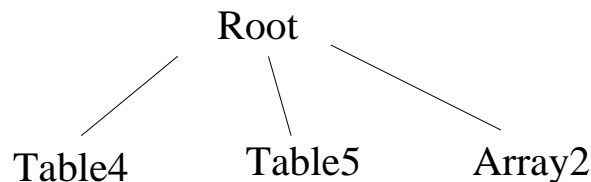


Figure 5.8: Resulting object tree derived from the use of the `rootUEP` parameter.

5.4 Selecting an User Entry Point (UEP) in your tree

If you have a **huge** tree in your data file with many nodes on it, creating the object tree would take long time. Many times, however, you are interested only in access to a part of the complete tree, so you won't strictly need PyTables to build the entire object tree in-memory, but only the *interesting* part.

This is where the `rootUEP` parameter of `openFile()` function (see 4.1.2) can be helpful. Imagine that you have a file called `"test.h5"` with the associated tree that you can see in figure 5.7, and you are interested only in the section marked in red. You can avoid the build of all the object tree by saying to `openFile` that your root will be the `/Group2/Group3` group. That is:

```
fileh = openFile("test.h5", rootUEP="/Group2/Group3")
```

As a result, the actual object tree built will be like the one that can be seen in figure 5.8.

Of course this has been a simple example and the use of the `rootUEP` parameter was not very necessary. But when you have *thousands* of nodes on a tree, you will certainly appreciate the `rootUEP` parameter.

Appendix A

Supported data types in tables

`IsDescription` subclasses supports a limited set of data types to define the table fields. Such a set is roughly the same than the types supported by the `numpy` package (see Greenfield *et al.*) in Python, with the exception of the complex datatypes that are not supported yet.

These data types in table columns can be set through the use of the `Col` class and its descendants (see 4.3). You may find useful the table A as a quick reference to the complete set of supported data types in PyTables.

Table A.1: Data types supported by subclasses of `IsDescription` definitions.

Type Code	Description	C Type	Size (in bytes)	Python Counterpart
Int8	8-bit integer	signed char	1	Integer
UInt8	8-bit unsigned integer	unsigned char	1	Integer
Int16	16-bit integer	short	2	Integer
UInt16	16-bit unsigned integer	unsigned short	2	Integer
Int32	integer	int	4	Integer
UInt32	unsigned integer	unsigned int	4	Long
Int64	64-bit integer	long long	8	Long
UInt64	unsigned 64-bit integer	unsigned long long	8	Long
Float32	single-precision float	float	4	Float
Float64	double-precision float	double	8	Float
CharType	arbitrary length string	char[]	*	String

Bibliography

- ASCHER, David, Paul F. DUBOIS, Konrad HINSEN, Jim HUGUNIN, and Travis OLIPHANT, : *Numerical Python*. Package to speed-up arithmetic operations on arrays of numbers.
URL <http://www.pfdubois.com/numpy/> 2, 7
- EWING, Greg, : *Pyrex. A Language for Writing Python Extension Modules*.
URL <http://www.cosc.canterbury.ac.nz/~greg/python/Pyrex> 7
- GAILLY, JeanLoup and Mark ADLER, : *zlib. A Massively Spiffy Yet Delicately Unobtrusive Compression Library*. A standard library for compression purposes.
URL <http://www.gzip.org/zlib/> 47
- GREENFIELD, Perry, Todd MILLER, Richard L. WHITE, and J. C. HSU., : *Numarray*. Reimplementation of Numeric which adds the ability to efficiently manipulate large numeric arrays in ways similar to Matlab and IDL. Among others, Numarray provides the record array extension.
URL <http://stsdas.stsci.edu/numarray/> 2, 7, 42, 51
- MERTZ, David, : *Objectify. On the 'Pythonic' treatment of XML documents as objects(II)*. Article describing XML Objectify, a Python module that allows working with XML documents as Python objects. Some of the ideas presented here are used in PyTables.
URL <http://www-106.ibm.com/developerworks/xml/library/xml-matters2/index.html> 2
- NCSA, : *What is HDF5?* Concise description about HDF5 capabilities and its differences from earlier versions (HDF4).
URL <http://hdf.ncsa.uiuc.edu/whatishdf5.html> 1
- OBERHUMER, Markus F.X.J., : *LZO and UCL. A couple of portable lossless data compression libraries*. They offer pretty fast compression and extremely fast decompression.
URL <http://www.oberhumer.com/opensource/> 7, 47
- RIGO, Armin, : *Psyco. A Python specializing compiler*. Run existing Python software faster, with no change in your source.
URL <http://psyco.sourceforge.net> 45