

Hierarchical datasets in Python

PyTables User Guide

Release 3.5.2

PyTables maintainers

October 22, 2019

CONTENTS

1	The PyTables Core Library	7
2	Complementary modules	101
3	Appendixes	107
	Index	127

LIST OF FIGURES

LIST OF TABLES

Authors Francesc Alted, Ivan Vilata, Scott Prater, Vicent Mas, Tom Hedley, Antonio Valentino, Jeffrey Whitaker, Anthony Scopatz, Josh Moore

Copyright © 2002, 2003, 2004 - Francesc Alted

© 2005, 2006, 2007 - Cárabos Coop. V.

© 2008, 2009, 2010 - Francesc Alted

© 2011-2018 - PyTables maintainers

Date October 22, 2019

Version 3.5.2

Home Page <http://www.pytables.org>

Copyright Notice and Statement for PyTables User's Guide

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- a. Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- b. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- c. Neither the name of Francesc Alted nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "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 THE COPYRIGHT OWNER OR CONTRIBUTORS 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.

THE PYTABLES CORE LIBRARY

1.1 Introduction

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

[Wisdom isn't worth anything if you can't use it to come up with a new way to cook garbanzos.]

—Gabriel García Márquez, A wise Catalan in “*Cien años de soledad*”

The goal of PyTables is to enable the end user to manipulate easily data *tables* and *array* objects in a hierarchical structure. The foundation of the underlying hierarchical data organization is the excellent HDF5 library (see [\[HDGF1\]](#)).

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

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* may seem to be 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 data acquisition systems, Internet services or scientific applications, for example) in an efficient manner that reduces demand on CPU time and I/O.

In order to emulate in Python records mapped to HDF5 C structs PyTables implements a special class so as to easily define all its fields and other properties. PyTables also provides a powerful interface to mine data in tables. 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 named fields and type information, such as in the following example:

```
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    = Int32Col()         # integer

    # A sub-structure (nested data-type)
    class Properties(IsDescription):
        pressure = Float32Col(shape=(2,3)) # 2-D float array (single-precision)
        energy   = Float64Col(shape=(2,3,4)) # 3-D float array (double-precision)
```

You then pass this class to the table constructor, fill its rows with your values, and save (arbitrarily large) collections of them to a file for persistent storage. After that, the 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 a library to interface with HDF5).

Other important entities in PyTables are *array* objects, which are analogous to tables with the difference that all of their components are homogeneous. They come in different flavors, like *generic* (they provide a quick and fast way to deal with for numerical arrays), *enlargeable* (arrays can be extended along a single dimension) and *variable length* (each row in the array can have a different number of elements).

The next section describes the most interesting capabilities of PyTables.

1.1.1 Main Features

PyTables takes advantage of the object orientation and introspection capabilities offered by Python, the powerful data management features of HDF5, and NumPy's flexibility and Numexpr's high-performance manipulation of large sets of objects organized in a grid-like fashion to provide these features:

- *Support for table entities:* You can tailor your data adding or deleting records in your tables. Large numbers of rows (up to 2^{63} , much more than will fit into memory) are supported as well.
- *Multidimensional and nested table cells:* You can declare a column to consist of values having any number of dimensions besides scalars, which is the only dimensionality allowed by the majority of relational databases. You can even declare columns that are made of other columns (of different types).
- *Indexing support for columns of tables:* Very useful if you have large tables and you want to quickly look up for values in columns satisfying some criteria.
- *Support for numerical arrays:* NumPy (see [\[NUMPY\]](#)) arrays can be used as a useful complement of tables to store homogeneous data.
- *Enlargeable arrays:* You can add new elements to existing arrays on disk in any dimension you want (but only one). Besides, you are able to access just a slice of your datasets by using the powerful extended slicing mechanism, without need to load all your complete dataset in memory.
- *Variable length arrays:* The number of elements in these arrays can vary from row to row. This provides a lot of flexibility when dealing with complex data.
- *Supports a hierarchical data model:* Allows the user to clearly structure all data. PyTables builds up an *object tree* in memory that replicates the underlying file data structure. Access to objects in the file is achieved by walking through and manipulating this object tree. Besides, this object tree is built in a lazy way, for efficiency purposes.
- *User defined metadata:* Besides supporting system metadata (like the number of rows of a table, shape, flavor, etc.) the user may specify arbitrary metadata (as for example, room temperature, or protocol for IP traffic that was collected) that complement the meaning of actual data.
- *Ability to read/modify generic HDF5 files:* PyTables can access a wide range of objects in generic HDF5 files, like compound type datasets (that can be mapped to Table objects), homogeneous datasets (that can be mapped to Array objects) or variable length record datasets (that can be mapped to VArray objects). Besides, if a dataset is not supported, it will be mapped to a special *UnImplemented* class (see [The UnImplemented class](#)), that will let the user see that the data is there, although it will be unreachable (still, you will be able to access the attributes and some metadata in the dataset). With that, PyTables probably can access and *modify* most of the HDF5 files out there.
- *Data compression:* Supports data compression (using the *Zlib*, *LZO*, *bzip2* and *Blosc* compression libraries) out of the box. This is important when you have repetitive data patterns and don't want to spend time searching for an optimized way to store them (saving you time spent analyzing your data organization).
- *High performance I/O:* On modern systems storing 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 that limit is surmountable!

- *Support of files bigger than 2 GB:* PyTables automatically inherits this capability from the underlying HDF5 library (assuming your platform supports the C long long integer, or, on Windows, `__int64`).
- *Architecture-independent:* PyTables has been carefully coded (as HDF5 itself) with little-endian/big-endian byte ordering issues in mind. So, you can write a file on a big-endian machine (like a Sparc or MIPS) and read it on other little-endian machine (like an Intel or Alpha) without problems. In addition, it has been tested successfully with 64 bit platforms (Intel-64, AMD-64, PowerPC-G5, MIPS, UltraSparc) using code generated with 64 bit aware compilers.

1.1.2 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. The HDF5 objects are read by walking through this object tree. You can get a good picture of what kind of data is kept in the object by examining the *metadata* nodes.

The different nodes in the object tree are instances of PyTables classes. There are several types of classes, but the most important ones are the Node, Group and Leaf classes. All nodes in a PyTables tree are instances of the Node class. The Group and Leaf classes are descendants of Node. Group instances (referred to as *groups* from now on) are a grouping structure containing instances of zero or more groups or leaves, together with supplementary metadata. Leaf instances (referred to as *leaves*) are containers for actual data and can not contain further groups or leaves. The Table, Array, CArray, EArray, VArray and UnImplemented classes are descendants of Leaf, and inherit all its properties.

Working with groups and leaves is similar in many ways to working with directories and files on a Unix filesystem, i.e. a node (file or directory) is always a *child* of one and only one group (directory), its *parent group*¹. Inside of that group, the node is accessed by its *name*. As is the case with Unix directories and files, objects in the object tree are often referenced by giving their full (absolute) path names. In PyTables this full path can be specified either as string (such as `"/subgroup2/table3"`, using `/` as a parent/child separator) or as a complete object path written in a format known as the *natural name* schema (such as `file.root.subgroup2.table3`).

Support for *natural naming* is a key aspect of PyTables. It means that the names of instance variables of the node objects are the same as the names of its children². This is very *Pythonic* and intuitive in many cases. Check the tutorial *Reading (and selecting) data in a table* for usage examples.

You should also be aware that not all the data present in a file is loaded into the object tree. The *metadata* (i.e. special data that describes the structure of the actual data) is loaded only when the user want to access to it (see later). Moreover, the actual data is not read until she request it (by calling a method on a particular node). Using the object tree (the metadata) you can retrieve information about the objects on disk such as table names, titles, column names, data types in columns, numbers of rows, or, in the case of arrays, their shapes, typecodes, etc. You can also search through the tree for specific kinds of data then read it and process it. In a certain sense, you can think of PyTables as a tool that applies the same introspection capabilities of Python objects to large amounts of data in persistent storage.

It is worth noting that PyTables sports a *metadata cache system* that loads nodes *lazily* (i.e. on-demand), and unloads nodes that have not been used for some time (following a *Least Recently Used* schema). It is important to stress out that the nodes enter the cache after they have been unreferenced (in the sense of Python reference counting), and that they can be revived (by referencing them again) directly from the cache without performing the de-serialization process from disk. This feature allows dealing with files with large hierarchies very quickly and with low memory consumption, while retaining all the powerful browsing capabilities of the previous implementation of the object tree. See [\[OPTIM\]](#) for more facts about the advantages introduced by this new metadata cache system.

To better understand the dynamic nature of this object tree entity, let's start with a sample PyTables script (which you can find in `examples/objecttree.py`) to create an HDF5 file:

¹ PyTables does not support hard links - for the moment.

² I got this simple but powerful idea from the excellent Objectify module by David Mertz (see [\[MERTZ\]](#)).

```
from tables import *

class Particle(IsDescription):
    identity = StringCol(itemsize=22, dflt=" ", pos=0) # character String
    idnumber = Int16Col(dflt=1, pos = 1) # short integer
    speed    = Float32Col(dflt=1, pos = 2) # single-precision

# Open a file in "w"rite mode
fileh = open_file("objecttree.h5", mode = "w")

# Get the HDF5 root group
root = fileh.root

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

# Now, create an array in root group
array1 = fileh.create_array(root, "array1", ["string", "array"], "String array")

# Create 2 new tables in group1
table1 = fileh.create_table(group1, "table1", Particle)
table2 = fileh.create_table("/group2", "table2", Particle)

# Create the last table in group2
array2 = fileh.create_array("/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 in [Figure 1³](#). When the file is created, the metadata in the object tree is updated in memory while the actual data is saved to disk. When you close the file the object tree is no longer available. However, when you reopen this file the object tree will be reconstructed in memory from the metadata on disk (this is done in a lazy way, in order to load only the objects that are required by the user), allowing you to work with it in exactly the same way as when you originally created it.

In [Figure 2](#), you can see an example of the object tree created when the above `objecttree.h5` file is read (in fact, such an object tree is always created when reading any supported generic HDF5 file). It is worthwhile to take your time to

³ We have used ViTables (see [\[VITABLES\]](#)) in order to create this snapshot.

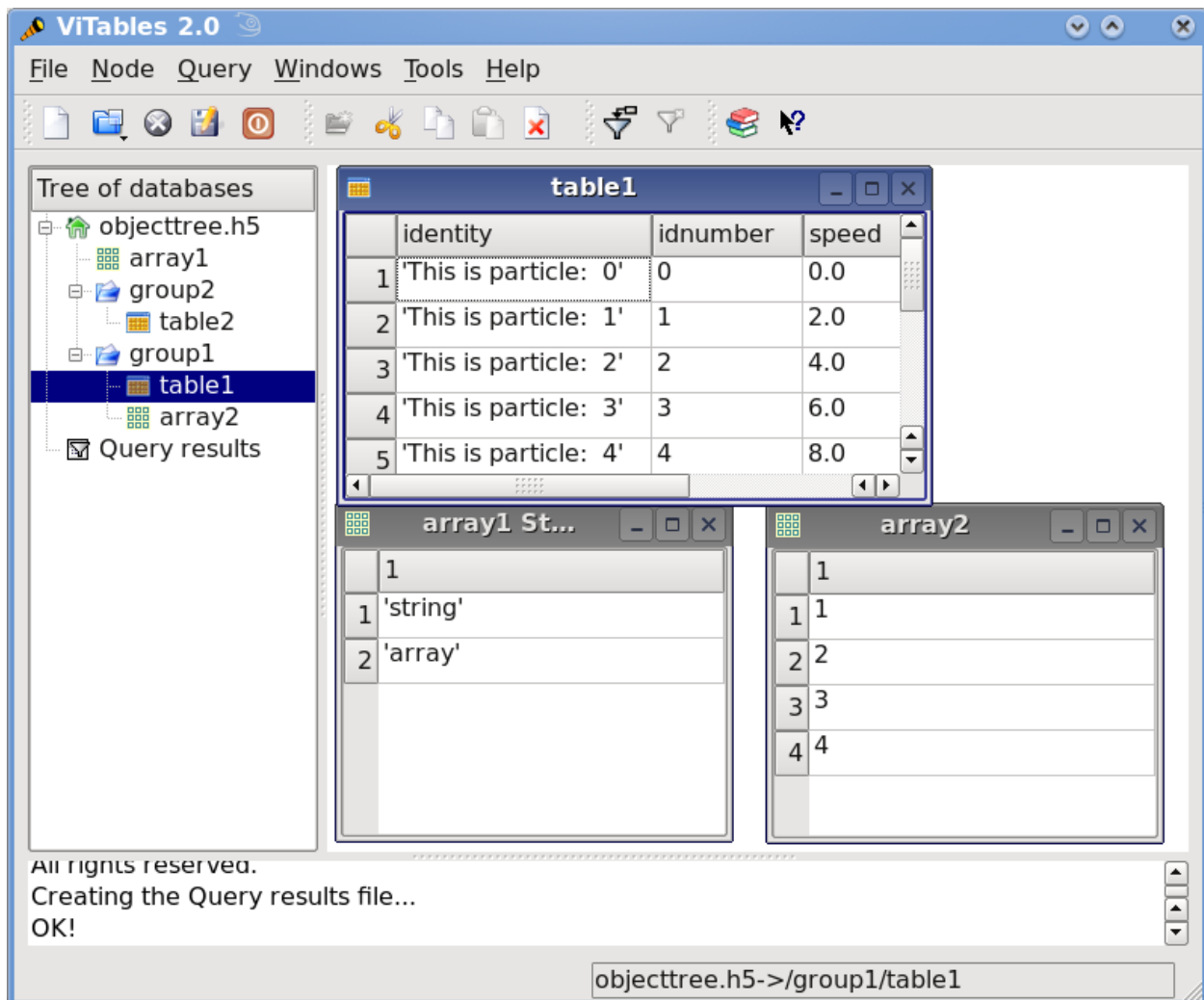


Fig. 1: Figure 1: An HDF5 example with 2 subgroups, 2 tables and 1 array.

understand it⁴. It will help you understand the relationships of in-memory PyTables objects.

1.2 Installation

Make things as simple as possible, but not any simpler.

—Albert Einstein

The Python Distutils are used to build and install PyTables, so it is fairly simple to get the application up and running. If you want to install the package from sources you can go on reading to the next section.

However, if you want to go straight to binaries that ‘just work’ for the main platforms (Linux, Mac OSX and Windows), you might want to use the excellent [Anaconda](#), [ActivePython](#), [Canopy](#) distributions. PyTables usually distributes its own Windows binaries too; go [Binary installation \(Windows\)](#) for instructions. Finally [Christoph Gohlke](#) also maintains an excellent suite of a variety of binary packages for Windows at his site.

1.2.1 Installation from source

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

Extensions in PyTables have been developed in Cython (see [\[CYTHON\]](#)) and the C language. You can rebuild everything from scratch if you have Cython installed, but this is not necessary, as the Cython compiled source is included in the source distribution.

To compile PyTables you will need a recent version of Python, the HDF5 (C flavor) library from <http://www.hdfgroup.org>, and the NumPy (see [\[NUMPY\]](#)) and Numexpr (see [\[NUMEXPR\]](#)) packages.

Prerequisites

First, make sure that you have

- [Python](#) ≥ 2.7 including Python 3.x
- [HDF5](#) $\geq 1.8.4$ ($\geq 1.8.15$ is strongly recommended)
- [NumPy](#) $\geq 1.9.3$
- [Numexpr](#) $\geq 2.6.2$
- [Cython](#) ≥ 0.21
- [c-blosc](#) $\geq 1.4.1$ (sources are bundled with PyTables sources but the user can use an external version of sources using the `BLOSC_DIR` environment variable or the `-blosc` flag of the `setup.py`)

installed (for testing purposes, we are using [HDF5 1.8.15](#), [NumPy 1.10.2](#) and [Numexpr 2.5.2](#) currently). If you don’t, fetch and install them before proceeding.

Compile and install these packages (but see [Windows prerequisites](#) for instructions on how to install pre-compiled binaries if you are not willing to compile the prerequisites on Windows systems).

⁴ 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.

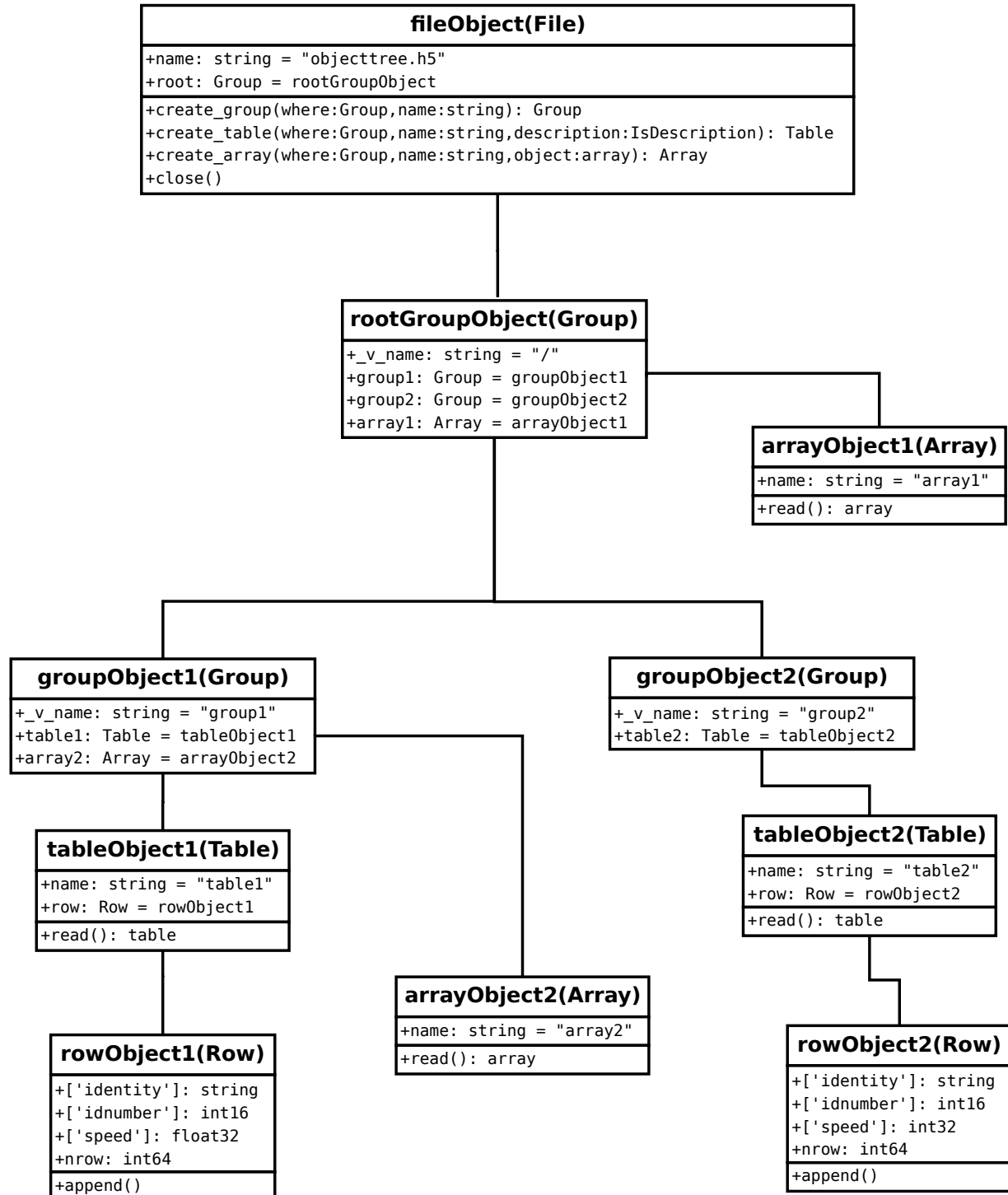


Fig. 2: Figure 2: A PyTables object tree example.

For compression (and possibly improved performance), you will need to install the Zlib (see [\[ZLIB\]](#)), which is also required by HDF5 as well. You may also optionally install the excellent LZO compression library (see [\[LZO\]](#) and [Compression issues](#)). The high-performance bzip2 compression library can also be used with PyTables (see [\[BZIP2\]](#)).

The Blosc (see [\[BLOSC\]](#)) compression library is embedded in PyTables, so this will be used in case it is not found in the system. So, in case the installer warns about not finding it, do not worry too much ;)

Unix

setup.py will detect HDF5, Blosc, LZO, or bzip2 libraries and include files under `/usr` or `/usr/local`; this will cover most manual installations as well as installations from packages. If setup.py can not find `libhdf5`, `libhdf5` (or `liblzo`, or `libbz2` that you may wish to use) or if you have several versions of a library installed and want to use a particular one, then you can set the path to the resource in the environment, by setting the values of the `HDF5_DIR`, `LZO_DIR`, `BZIP2_DIR` or `BLOSC_DIR` environment variables to the path to the particular resource. You may also specify the locations of the resource root directories on the setup.py command line. For example:

```
--hdf5=/stuff/hdf5-1.8.12
--blosc=/stuff/blosc-1.8.1
--lzo=/stuff/lzo-2.02
--bzip2=/stuff/bzip2-1.0.5
```

If your HDF5 library was built as a shared library 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/hdf5-1.8.12/lib"
```

You may also want to try setting the `LD_LIBRARY_PATH` environment variable to point to the directory where the shared libraries can be found. Check your compiler and linker documentation as well as the Python Distutils documentation for the correct syntax or environment variable names. It is also possible to link with specific libraries by setting the `LIBS` environment variable:

```
LIBS="hdf5-1.8.12 ns1"
```

Starting from PyTables 3.2 can also query the `pkg-config` database to find the required packages. If available, `pkg-config` is used by default unless explicitly disabled.

To suppress the use of `pkg-config`:

```
$ python setup.py build --use-pkgconfig=FALSE
```

or use the `USE_PKGCONFIG` environment variable:

```
$ env USE_PKGCONFIG=FALSE python setup.py build
```

Windows

You can get ready-to-use Windows binaries and other development files for most of the following libraries from the GnuWin32 project (see [\[GNUWIN32\]](#)). In case you cannot find the LZO binaries in the GnuWin32 repository, you can find them at <http://sourceforge.net/projects/pytables/files/lzo-win>. Once you have installed the prerequisites, setup.py needs to know where the necessary library *stub* (.lib) and *header* (.h) files are installed. You can set the path to the include and dll directories for the HDF5 (mandatory) and LZO, BZIP2, BLOSC (optional) libraries in the environment, by setting the values of the `HDF5_DIR`, `LZO_DIR`, `BZIP2_DIR` or `BLOSC_DIR` environment variables to the path to the particular resource. For example:

```
set HDF5_DIR=c:\\stuff\\hdf5-1.8.5-32bit-VS2008-IVF101\\release
set BLOSC_DIR=c:\\Program Files (x86)\\Blosc
```

(continues on next page)

(continued from previous page)

```
set LZO_DIR=c:\\Program Files (x86)\\GnuWin32
set BZIP2_DIR=c:\\Program Files (x86)\\GnuWin32
```

You may also specify the locations of the resource root directories on the setup.py command line. For example:

```
--hdf5=c:\\stuff\\hdf5-1.8.5-32bit-VS2008-IVF101\\release
--blosc=c:\\Program Files (x86)\\Blosc
--lzo=c:\\Program Files (x86)\\GnuWin32
--bzip2=c:\\Program Files (x86)\\GnuWin32
```

Conda

Pre-built packages for PyTables are available in the anaconda (default) channel:

```
conda install pytables
```

The most recent version is usually available in the conda-forge channel:

```
conda config --add channels conda-forge
conda install pytables
```

The HDF5 libraries and other helper packages are automatically found in a conda environment. During installation setup.py uses the `CONDA_PREFIX` environment variable to detect a conda environment. If detected it will try to find all packages within this environment. PyTables needs at least the hdf5 package:

```
conda install hdf5
python setup.py install
```

It is still possible to override package locations using the `HDF5_DIR`, `LZO_DIR`, `BZIP2_DIR` or `BLOSC_DIR` environment variables.

When inside a conda environment `pkg-config` will not work. To disable using the conda environment and fall back to `pkg-config` use `--no-conda`:

```
python setup.py install --no-conda
```

When the `--use-pkgconfig` flag is used, `--no-conda` is assumed.

Development version (Unix)

Installation of the development version is very similar to installation from a source package (described above). There are two main differences:

1. sources have to be downloaded from the [PyTables source repository](#) hosted on [GitHub](#). Git (see [\[GIT\]](#)) is used as VCS. The following command create a local copy of latest development version sources:

```
$ git clone https://github.com/PyTables/PyTables.git
```

2. sources in the git repository do not include pre-built documentation and pre-generated C code of Cython extension modules. To be able to generate them, both Cython (see [\[CYTHON\]](#)) and sphinx $\geq 1.0.7$ (see [\[SPHINX\]](#)) are mandatory prerequisites.

PyTables package installation

Once you have installed the HDF5 library and the NumPy and Numexpr packages, you can proceed with the PyTables package itself.

1. Run this command from the main PyTables distribution directory, including any extra command line arguments as discussed above:

```
$ python setup.py build
```

If the HDF5 installation is in a custom path, e.g. \$HOME/hdf5-1.8.15pre7, one of the following commands can be used:

```
$ python setup.py build --hdf5=$HOME/hdf5-1.8.15pre7
```

Note: AVX2 support is detected automatically for your machine and, if found, it is enabled by default. In some situations you may want to disable AVX2 explicitly (maybe your binaries have to be exported and run on machines that do not have AVX2 support). In that case, define the `DISABLE_AVX2` environment variable:

```
$ DISABLE_AVX2=True python setup.py build # for bash and its variants
```

2. To run the test suite, execute any of these commands.

Unix In the sh shell and its variants:

```
$ cd build/lib.linux-x86_64-3.3
$ env PYTHONPATH=. python tables/tests/test_all.py
```

or, if you prefer:

```
$ cd build/lib.linux-x86_64-3.3
$ env PYTHONPATH=. python -c "import tables; tables.test()"
```

Note: the syntax used above overrides original contents of the `PYTHONPATH` environment variable. If this is not the desired behaviour and the user just wants to add some path before existing ones, then the safest syntax to use is the following:

```
$ env PYTHONPATH=.${PYTHONPATH:+:$PYTHONPATH} python tables/tests/test_all.py
```

Please refer to your **sh** documentation for details.

Windows

Open the command prompt (cmd.exe or command.com) and type:

```
> cd build\\lib.linux-x86_64-2.7
> set PYTHONPATH=.;%PYTHONPATH%
> python tables\\tests\\test_all.py
```

or:

```
> cd build\\lib.linux-x86_64-2.7
> set PYTHONPATH=.;%PYTHONPATH%
> python -c "import tables; tables.test()"
```

Both commands do the same thing, but the latter still works on an already installed PyTables (so, there is no need to set the `PYTHONPATH` variable for this case). However, before installation, the former is recommended because it is more flexible, as you can see below. If you would like to see verbose output from the tests simply add the `-v` flag and/or the word `verbose` to the first of the command lines above. You can also run only the tests in a particular test module. For example, to execute just the `test_types` test suite, you only have to specify it:

```
# change to backslashes for win
$ python tables/tests/test_types.py -v
```

You have other options to pass to the `test_all.py` driver:

```
# change to backslashes for win
$ python tables/tests/test_all.py --heavy
```

The command above runs every test in the test unit. Beware, it can take a lot of time, CPU and memory resources to complete:

```
# change to backslashes for win
$ python tables/tests/test_all.py --print-versions
```

The command above shows the versions for all the packages that PyTables relies on. Please be sure to include this when reporting bugs:

```
# only under Linux 2.6.x
$ python tables/tests/test_all.py --show-memory
```

The command above prints out the evolution of the memory consumption after each test module completion. It's useful for locating memory leaks in PyTables (or packages behind it). Only valid for Linux 2.6.x kernels. And last, but not least, in case a test fails, please run the failing test module again and enable the verbose output:

```
$ python tables/tests/test_<module>.py -v verbose
```

and, very important, obtain your PyTables version information by using the `--print-versions` flag (see above) and send back both outputs to developers so that we may continue improving PyTables. If you run into problems because Python can not load the HDF5 library or other shared libraries.

Unix

Try setting the `LD_LIBRARY_PATH` or equivalent environment variable to point to the directory where the missing libraries can be found.

Windows

Put the DLL libraries (`hdf5dll.dll` and, optionally, `lzo1.dll`, `bzip2.dll` or `blosc.dll`) in a directory listed in your `PATH` environment variable. The `setup.py` installation program will print out a warning to that effect if the libraries can not be found.

3. To install the entire PyTables Python package, change back to the root distribution directory and run the following command (make sure you have sufficient permissions to write to the directories where the PyTables files will be installed):

```
$ python setup.py install
```

Again if one needs to point to libraries installed in custom paths, then specific `setup.py` options can be used:

```
$ python setup.py install --hdf5=/hdf5/custom/path
```

or:

```
$ env HDF5_DIR=/hdf5/custom/path python setup.py install
```

Of course, you will need super-user privileges if you want to install PyTables on a system-protected area. You can select, though, a different place to install the package using the *-prefix* flag:

```
$ python setup.py install --prefix="/home/myuser/mystuff"
```

Have in mind, however, that if you use the *-prefix* flag to install in a non-standard place, you should properly setup your `PYTHONPATH` environment variable, so that the Python interpreter would be able to find your new PyTables installation. You have more installation options available in the Distutils package. Issue a:

```
$ python setup.py install --help
```

for more information on that subject.

That's it! Now you can skip to the next chapter to learn how to use PyTables.

1.2.2 Installation with **pip**

Many users find it useful to use the **pip** program (or similar ones) to install python packages.

As explained in previous sections the user should in any case ensure that all dependencies listed in the *Prerequisites* section are correctly installed.

The simplest way to install PyTables using **pip** is the following:

```
$ pip install tables
```

The following example shows how to install the latest stable version of PyTables in the user folder when a older version of the package is already installed at system level:

```
$ pip install --user --upgrade tables
```

The *-user* option tells to the **pip** tool to install the package in the user folder (`$HOME/.local` on GNU/Linux and Unix systems), while the *-upgrade* option forces the installation of the latest version even if an older version of the package is already installed.

Additional options for the `setup.py` script can be specified using them *-install-option*:

```
$ pip install --install-option='--hdf5=/custom/path/to/hdf5' tables
```

or:

```
$ env HDF5_DIR=/custom/path/to/hdf5 pip install tables
```

The **pip** tool can also be used to install packages from a source tar-ball:

```
$ pip install tables-3.0.0.tar.gz
```

To install the development version of PyTables from the *develop* branch of the main **git** [\[GIT\]](#) repository the command is the following:

```
$ pip install git+https://github.com/PyTables/PyTables.git@develop#egg=tables
```

A similar command can be used to install a specific tagged fersion:

```
$ pip install git+https://github.com/PyTables/PyTables.git@v.2.4.0#egg=tables
```

Finally, PyTables developers provide a `requirements.txt` file that can be used by **pip** to install the PyTables dependencies:

```
$ wget https://raw.githubusercontent.com/PyTables/PyTables/develop/requirements.txt
$ pip install -r requirements.txt
```

Of course the `requirements.txt` file can be used to install only python packages. Other dependencies like the HDF5 library or compression libraries have to be installed by the user.

Note: Recent versions of [Debian](#) and [Ubuntu](#) the HDF5 library is installed in with a very peculiar layout that allows to have both the serial and MPI versions installed at the same time.

PyTables >= 3.2 natively supports the new layout via `pkg-config` (that is expected to be installed on the system at build time).

If `pkg-config` is not available or PyTables is older than version 3.2, then the following command can be used:

```
$ env CPPFLAGS=-I/usr/include/hdf5/serial \
LDFLAGS=-L/usr/lib/x86_64-linux-gnu/hdf5/serial python3 setup.py install
```

or:

```
$ env CPPFLAGS=-I/usr/include/hdf5/serial \
LDFLAGS=-L/usr/lib/x86_64-linux-gnu/hdf5/serial pip install tables
```

1.2.3 Binary installation (Windows)

This section is intended for installing precompiled binaries on Windows platforms. Binaries are distributed in wheel format, which can be downloaded and installed using pip as described above. You may also find it useful for instructions on how to install *binary prerequisites* even if you want to compile PyTables itself on Windows.

Windows prerequisites

First, make sure that you have Python 2.7, NumPy 1.8.0 and Numexpr 2.5.2 or higher installed.

To enable compression with the optional LZO library (see the [Compression issues](#) for hints about how it may be used to improve performance), fetch and install the LZO from <http://sourceforge.net/projects/pytables/files/lzo-win> (choose v1.x for Windows 32-bit and v2.x for Windows 64-bit). Normally, you will only need to fetch that package and copy the included `lzo1.dll/lzo2.dll` file in a directory in the `PATH` environment variable (for example `C:\WINDOWS\SYSTEM`) or `python_installation_path\Lib\site-packages\tables` (the last directory may not exist yet, so if you want to install the DLL there, you should do so *after* installing the PyTables package), so that it can be found by the PyTables extensions.

Please note that PyTables has internal machinery for dealing with uninstalled optional compression libraries, so, you don't need to install the LZO or bzip2 dynamic libraries if you don't want to.

PyTables package installation

On PyPI wheels for 32 and 64-bit versions of Windows are usually provided. They are automatically found and installed using pip:

```
$ pip install tables
```

If a matching wheel cannot be found for your installation, third party built wheels can be found e.g. at the [Unofficial Windows Binaries for Python Extension Packages](#) page. Download the wheel matching the version of python and either the 32 or 64-bit version and install using pip:

```
# python 3.5 64-bit:
$ pip install tables-3.3-cp35-cp35m-win_amd64.whl
```

You can (and *you should*) test your installation by running the next commands:

```
>>> import tables
>>> tables.test()
```

on your favorite python shell. If all the tests pass (possibly with a few warnings, related to the potential unavailability of LZO lib) you already have a working, well-tested copy of PyTables installed! If any test fails, please copy the output of the error messages as well as the output of:

```
>>> tables.print_versions()
```

and mail them to the developers so that the problem can be fixed in future releases.

You can proceed now to the next chapter to see how to use PyTables.

1.3 Tutorials

Seràs la clau que obre tots els panys, seràs la llum, la llum il.limitada, seràs confí on l'aurora comença,
seràs forment, escala il.luminada!

—Lyrics: Vicent Andrés i Estellés. Music: Ovidi Montllor, Toti Soler, M'aclame a tu

This chapter consists of a series of simple yet comprehensive tutorials that will enable you to understand PyTables' main features. If you would like more information about some particular instance variable, global function, or method, look at the doc strings or go to the library reference in [Library Reference](#). If you are reading this in PDF or HTML formats, follow the corresponding hyperlink near each newly introduced entity.

Please note that throughout this document the terms *column* and *field* will be used interchangeably, as will the terms *row* and *record*.

1.3.1 Getting started

In this section, we will see how to define our own records in Python and save collections of them (i.e. a *table*) into a file. Then we will select some of the data in the table using Python cuts and create NumPy arrays to store this selection as separate objects in a 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. I encourage you to do parallel testing and inspect the created objects (variables, docs, children objects, etc.) during the course of the tutorial!

Importing tables objects

Before starting you need to import the public objects in the tables package. You normally do that by executing:

```
>>> import tables
```

This is the recommended way to import tables if you don't want to pollute your namespace. However, PyTables has a contained set of first-level primitives, so you may consider using the alternative:

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

If you are going to work with NumPy arrays (and normally, you will) you will also need to import functions from the numpy package. So most PyTables programs begin with:

```
>>> import tables      # but in this tutorial we use "from tables import *"
>>> import numpy
```

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 retrieved from it. You need first to define the table, the number of columns it has, what kind of object is contained in each column, and so on.

Our particle detector has a TDC (Time to Digital Converter) counter with a dynamic range of 8 bits and an ADC (Analogical 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, so we will add two new fields called grid_i and grid_j. Our instrumentation also can obtain the pressure and energy of the particle. The resolution of the pressure-gauge allows us to use a single-precision float to store pressure readings, while the energy value will need a double-precision float. Finally, to track the particle we want to assign it a name to identify the kind of the particle it is and a unique numeric identifier. So we will add two more fields: name will be a string of up to 16 characters, and idnumber will be an integer of 64 bits (to allow us to store records for extremely large numbers of particles).

Having determined our columns and their types, we can now declare a new Particle class that will contain all this information:

```
>>> from tables import *
>>> 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()         # 32-bit integer
...     grid_j    = Int32Col()         # 32-bit integer
...     pressure  = Float32Col()       # float (single-precision)
...     energy    = Float64Col()       # double (double-precision)
>>>
```

This definition class is self-explanatory. Basically, you declare a class variable for each field you need. As its value you assign an instance of the appropriate Col subclass, according to the kind of column defined (the data type, the length, the shape, etc). See the [The Col class and its descendants](#) for a complete description of these subclasses. See also [Supported data types in PyTables](#) for a list of data types supported by the Col constructor.

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

Creating a PyTables file from scratch

Use the top-level `open_file()` function to create a PyTables file:

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

`open_file()` is one of the objects imported by the `from tables import *` statement. Here, we are saying that we want to create a new file in the current working directory called “tutorial1.h5” in “w”rite mode and with an descriptive title string (“Test file”). This function attempts to open the file, and if successful, returns the File (see [The File Class](#)) object instance `h5file`. The root of the object tree is specified in the instance’s `root` attribute.

Creating a new group

Now, to better organize our data, we will create a group called *detector* that branches from the root node. We will save our particle data table in this group:

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

Here, we have taken the File instance `h5file` and invoked its `File.create_group()` method to create a new group called *detector* branching from “/” (another way to refer to the `h5file.root` object we mentioned above). This will create a new Group (see [The Group class](#)) object instance that will be assigned to the variable `group`.

Creating a new table

Let’s now create a Table (see [The Table class](#)) object as a branch off the newly-created group. We do that by calling the `File.create_table()` method of the `h5file` object:

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

We create the Table instance under `group`. We assign this table the node name “*readout*”. The `Particle` class declared before is the *description* parameter (to define the columns of the table) and finally we set “*Readout example*” as the Table title. With all this information, a new Table instance is created and assigned to the variable `table`.

If you are curious about how the object tree looks right now, simply print the File instance variable `h5file`, and examine the output:

```
>>> print(h5file)
tutorial1.h5 (File) 'Test file'
Last modif.: 'Wed Mar  7 11:06:12 2007'
Object Tree:
/ (RootGroup) 'Test file'
/detector (Group) 'Detector information'
/detector/readout (Table(0,)) 'Readout example'
```

As you can see, a dump of the object tree is displayed. It’s easy to see the Group and Table objects we have just created. If you want more information, just type the variable containing the File instance:

```
>>> h5file
File(filename='tutorial1.h5', title='Test file', mode='w', root_uep='/',
filters=Filters(complevel=0, shuffle=False, bitshuffle=False, fletcher32=False))
/ (RootGroup) 'Test file'
/detector (Group) 'Detector information'
/detector/readout (Table(0,)) 'Readout example'
description := {
  "ADCcount": UInt16Col(shape=(), dflt=0, pos=0),
```

(continues on next page)

(continued from previous page)

```
"TDCcount": UInt8Col(shape=(), dflt=0, pos=1),
"energy": Float64Col(shape=(), dflt=0.0, pos=2),
"grid_i": Int32Col(shape=(), dflt=0, pos=3),
"grid_j": Int32Col(shape=(), dflt=0, pos=4),
"idnumber": Int64Col(shape=(), dflt=0, pos=5),
"name": StringCol(itemsize=16, shape=(), dflt='', pos=6),
"pressure": Float32Col(shape=(), dflt=0.0, pos=7)}
byteorder := 'little'
chunkshape := (87,)
```

More detailed information is displayed about each object in the tree. Note how `Particle`, our table descriptor class, is printed as part of the *readout* table description information. In general, you can obtain much more information about the objects and their children by just printing them. That introspection capability is very useful, and I recommend that you use it extensively.

The time has come to fill this table with some values. First we will get a pointer to the `Row` (see *The Row class*) instance of this table instance:

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

The `row` attribute of `table` points to the `Row` instance that will be used to write data rows into the table. We write data simply by assigning the `Row` instance the values for each row as if it were a dictionary (although it is actually an *extension class*), using the column names as keys.

Below is an example of how to write rows:

```
>>> 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)
...     # Insert a new particle record
...     particle.append()
>>>
```

This code should be easy to understand. The lines inside the loop just assign values to the different columns in the `Row` instance `particle` (see *The Row class*). A call to its `append()` method writes this information to the table I/O buffer.

After we have processed all our data, we should flush the table's I/O buffer if we want to write all this data to disk. We achieve that by calling the `table.flush()` method:

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

Remember, flushing a table is a *very important* step as it will not only help to maintain the integrity of your file, but also will free valuable memory resources (i.e. internal buffers) that your program may need for other things.

Reading (and selecting) data in a table

Ok. We have our data on disk, and now we need to access it and select from specific columns the values we are interested in. See the example below:

```
>>> 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 creates a “shortcut” to the *readout* table deeper on the object tree. As you can see, we use the *natural naming* schema to access it. We also could have used the `h5file.get_node()` method, as we will do later on.

You will recognize the last two lines as a Python list comprehension. It loops over the rows in *table* as they are provided by the `Table.iterrows()` iterator. The iterator returns values until all the data in *table* is exhausted. These rows are filtered using the expression:

```
x['TDCcount'] > 3 and 20 <= x['pressure'] < 50
```

So, we are selecting the values of the *pressure* column from filtered records to create the final list and assign it to *pressure* variable.

We could have used a normal for loop to accomplish the same purpose, but I find comprehension syntax to be more compact and elegant.

PyTables do offer other, more powerful ways of performing selections which may be more suitable if you have very large tables or if you need very high query speeds. They are called *in-kernel* and *indexed* queries, and you can use them through `Table.where()` and other related methods.

Let’s use an in-kernel selection to query the *name* column for the same set of cuts:

```
>>> names = [ x['name'] for x in table.where(""(TDCcount > 3) & (20 <= pressure) &
↳(pressure < 50) """) ]
>>> names
['Particle:      5', 'Particle:      6', 'Particle:      7']
```

In-kernel and indexed queries are not only much faster, but as you can see, they also look more compact, and are among the greatest features for PyTables, so be sure that you use them a lot. See [Condition Syntax](#) and [Accelerating your searches](#) for more information on in-kernel and indexed selections.

Note: A special care should be taken when the query condition includes string literals. Indeed Python 2 string literals are string of bytes while Python 3 strings are unicode objects.

With reference to the above definition of *Particle* it has to be noted that the type of the “name” column do not change depending on the Python version used (of course). It always corresponds to strings of bytes.

Any condition involving the “name” column should be written using the appropriate type for string literals in order to avoid `TypeError`s.

Suppose one wants to get rows corresponding to specific particle names.

The code below will work fine in Python 2 but will fail with a `TypeError` in Python 3:

```
>>> condition = '(name == "Particle:      5") | (name == "Particle:      7")'
>>> for record in table.where(condition): # TypeError in Python3
...     # do something with "record"
```

The reason is that in Python 3 “condition” implies a comparison between a string of bytes (“name” column contents) and an unicode literals.

The correct way to write the condition is:

```
>>> condition = '(name == b"Particle:      5") | (name == b"Particle:      7")'
```

That's enough about selections for now. The next section will show you how to save these selected results to a file.

Creating new array objects

In order to separate the selected data from the mass of detector data, we will create a new group columns branching off the root group. Afterwards, under this group, we will create two arrays that will contain the selected data. First, we create the group:

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

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

Now, create the first of the two Array objects we've just mentioned:

```
>>> h5file.create_array(gcolumns, 'pressure', array(pressure), "Pressure column_
↪selection")
/cOLUMNS/pressure (Array(3,)) 'Pressure column selection'
  atom := Float64Atom(shape=(), dflt=0.0)
  maInDim := 0
  flavor := 'numpy'
  byteorder := 'little'
  chunkshape := None
```

We already know the first two parameters of the `File.create_array()` methods (these are the same as the first two in `create_table`): they are the parent group *where* Array will be created and the Array instance *name*. The third parameter is the *object* we want to save to disk. In this case, it is a NumPy array that is built from the selection list we created before. The fourth parameter is the *title*.

Now, we will save the second array. It contains the list of strings we selected before: we save this object as-is, with no further conversion:

```
>>> h5file.create_array(gcolumns, 'name', names, "Name column selection")
/cOLUMNS/name (Array(3,)) 'Name column selection'
  atom := StringAtom(itemsizes=16, shape=(), dflt='')
  maInDim := 0
  flavor := 'python'
  byteorder := 'irrelevant'
  chunkshape := None
```

As you can see, `File.create_array()` accepts *names* (which is a regular Python list) as an *object* parameter. Actually, it accepts a variety of different regular objects (see `create_array()`) as parameters. The *flavor* attribute (see the output above) saves the original kind of object that was saved. Based on this *flavor*, PyTables will be able to retrieve exactly the same object from disk later on.

Note that in these examples, the `create_array` method returns an Array instance that is not assigned to any variable. Don't worry, this is intentional to show the kind of object we have created by displaying its representation. The Array objects have been attached to the object tree and saved to disk, as you can see if you print the complete object tree:

```
>>> print(h5file)
tutorial11.h5 (File) 'Test file'
Last modif.: 'Wed Mar  7 19:40:44 2007'
Object Tree:
```

(continues on next page)

(continued from previous page)

```

/ (RootGroup) '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'

```

Closing the file and looking at its content

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

```

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

```

You have now created your first PyTables file with a table and two arrays. You can examine it with any generic HDF5 tool, such as h5dump or h5ls. Here is what the tutorial1.h5 looks like when read with the h5ls program.

```

$ 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}

```

Here's the output as displayed by the “ptdump” PyTables utility (located in utils/ directory).

```

$ ptdump tutorial1.h5
/ (RootGroup) '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 ptdump if you want more verbosity. Try them out!

Also, in [Figure 1](#), you can admire how the tutorial1.h5 looks like using the ViTables graphical interface.

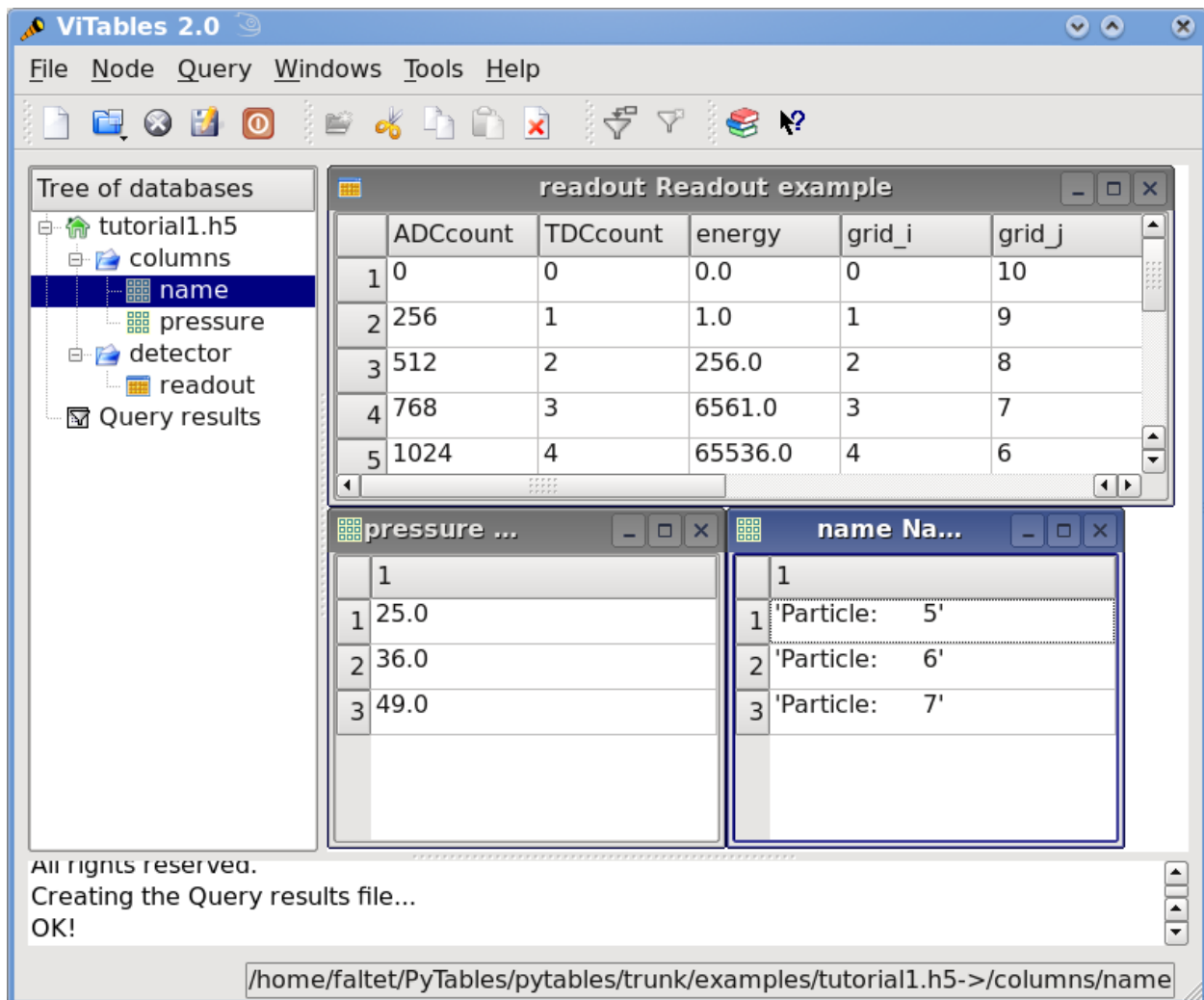


Fig. 3: Figure 1. The initial version of the data file for tutorial 1, with a view of the data objects.

1.3.2 Browsing the *object tree*

In this section, we will learn how to browse the tree and retrieve data and also meta-information about the actual data.

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 course of the tutorial.

Traversing the object tree

Let's start by opening the file we created in last tutorial section:

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

This time, we have opened the file in “a”ppend mode. We use this mode 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 search the tree.

To start with, you can get a preliminary overview of the object tree by simply printing the existing File instance:

```
>>> print(h5file)
tutorial1.h5 (File) 'Test file'
Last modif.: 'Wed Mar  7 19:50:57 2007'
Object Tree:
/ (RootGroup) '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'
```

It looks like all of 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)
/ (RootGroup) '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 `File.walk_groups()` method of the File class to list only the *groups* on tree:

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

Note that `File.walk_groups()` actually returns an *iterator*, not a list of objects. Using this iterator with the `list_nodes()` method is a powerful combination. Let's see an example listing of all the arrays in the tree:

```
>>> for group in h5file.walk_groups("/"):
...     for array in h5file.list_nodes(group, classname='Array'):
...         print(array)
```

(continues on next page)

(continued from previous page)

```
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
```

`File.list_nodes()` returns a list containing all the nodes hanging off a specific Group. If the *classname* keyword is specified, the method will filter out all instances which are not descendants of the class. We have asked for only Array instances. There exist also an iterator counterpart called `File.iter_nodes()` that might be handy in some situations, like for example when dealing with groups with a large number of nodes behind it.

We can combine both calls by using the `File.walk_nodes()` special method of the File object. For example:

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

This is a nice shortcut when working interactively.

Finally, we will list all the Leaf, i.e. Table and Array instances (see [The Leaf class](#) for detailed information on Leaf class), in the /detector group. Note that only one instance of the Table class (i.e. readout) will be selected in this group (as should be the case):

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

We have used a call to the `Group._f_walknodes()` method, using the *natural naming* path specification.

Of course you can do more sophisticated node selections using these powerful methods. But first, let's take a look at some important PyTables object instance variables.

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 [The AttributeSet class](#)). You can access 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:

```
>>> 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, let's 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 the `_v_attrs` attribute because `detector` is a Group node. In general, you can save any standard Python data structure as an attribute node. See [The AttributeSet class](#) for a more detailed explanation of how they are serialized for export to disk.

Retrieving the attributes is equally simple:

```
>>> 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 examine the current user attribute set of /detector/table, you can print its representation (try hitting the TAB key twice if you are on a Unix Python console with the rlcompleter module active):

```
>>> table.attrs
/detector/readout._v_attrs (AttributeSet), 23 attributes:
[CLASS := 'TABLE',
 FIELD_0_FILL := 0,
 FIELD_0_NAME := 'ADCcount',
 FIELD_1_FILL := 0,
 FIELD_1_NAME := 'TDCcount',
 FIELD_2_FILL := 0.0,
 FIELD_2_NAME := 'energy',
 FIELD_3_FILL := 0,
 FIELD_3_NAME := 'grid_i',
 FIELD_4_FILL := 0,
 FIELD_4_NAME := 'grid_j',
 FIELD_5_FILL := 0,
 FIELD_5_NAME := 'idnumber',
 FIELD_6_FILL := '',
 FIELD_6_NAME := 'name',
 FIELD_7_FILL := 0.0,
 FIELD_7_NAME := 'pressure',
 FLAVOR := 'numpy',
 NROWS := 10,
 TITLE := 'Readout example',
 VERSION := '2.6',
 temp_scale := 'Celsius',
 temperature := 18.399999999999999]
```

We've got all the attributes (including the *system* attributes). You can get a list of *all* attributes or only the *user* or *system* attributes with the `_f_list()` method:

```
>>> print(table.attrs._f_list("all"))
['CLASS', 'FIELD_0_FILL', 'FIELD_0_NAME', 'FIELD_1_FILL', 'FIELD_1_NAME',
'FIELD_2_FILL', 'FIELD_2_NAME', 'FIELD_3_FILL', 'FIELD_3_NAME', 'FIELD_4_FILL',
'FIELD_4_NAME', 'FIELD_5_FILL', 'FIELD_5_NAME', 'FIELD_6_FILL', 'FIELD_6_NAME',
'FIELD_7_FILL', 'FIELD_7_NAME', 'FLAVOR', 'NROWS', 'TITLE', 'VERSION',
'temp_scale', 'temperature']
>>> print(table.attrs._f_list("user"))
['temp_scale', 'temperature']
>>> print(table.attrs._f_list("sys"))
['CLASS', 'FIELD_0_FILL', 'FIELD_0_NAME', 'FIELD_1_FILL', 'FIELD_1_NAME',
'FIELD_2_FILL', 'FIELD_2_NAME', 'FIELD_3_FILL', 'FIELD_3_NAME', 'FIELD_4_FILL',
'FIELD_4_NAME', 'FIELD_5_FILL', 'FIELD_5_NAME', 'FIELD_6_FILL', 'FIELD_6_NAME',
```

(continues on next page)

(continued from previous page)

```
'FIELD_7_FILL', 'FIELD_7_NAME', 'FLAVOR', 'NROWS', 'TITLE', 'VERSION']
```

You can also rename attributes:

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

And, from PyTables 2.0 on, you are allowed also to set, delete or rename system attributes:

```
>>> table.attrs._f_rename("VERSION", "version")
>>> table.attrs.VERSION
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "tables/attributeset.py", line 222, in __getattr__
    (name, self._v__nodepath)
AttributeError: Attribute 'VERSION' does not exist in node: '/detector/readout'
>>> table.attrs.version
'2.6'
```

Caveat emptor: you must be careful when modifying system attributes because you may end fooling PyTables and ultimately getting unwanted behaviour. Use this only if you know what are you doing.

So, given the caveat above, we will proceed to restore the original name of VERSION attribute:

```
>>> table.attrs._f_rename("version", "VERSION")
>>> table.attrs.VERSION
'2.6'
```

Ok. that's better. If you would terminate your session now, you would be able to use the h5ls command to read the /detector/readout attributes from the file written to disk.

```
$ 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:      "2.6"
  Attribute: TITLE      scalar
    Type:      16-byte null-terminated ASCII string
    Data:      "Readout example"
  Attribute: NROWS      scalar
    Type:      native long long
    Data:      10
  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
```

(continues on next page)

(continued from previous page)

```

    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: FLAVOR scalar
    Type:      5-byte null-terminated ASCII string
    Data:      "numpy"
Attribute: tempScale scalar
    Type:      7-byte null-terminated ASCII string
    Data:      "Celsius"
Attribute: temperature scalar
    Type:      native double
    Data:      18.4
Location: 0:1:0:1952
Links: 1
Modified: 2006-12-11 10:35:13 CET
Chunks: {85} 3995 bytes
Storage: 470 logical bytes, 3995 allocated bytes, 11.76% 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

```

Attributes are a useful mechanism to add persistent (meta) information to your data.

Getting object metadata

Each object in PyTables has *metadata* information about the data in the file. Normally this *meta-information* is accessible through the node instance variables. Let's take 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:")

```

(continues on next page)

(continued from previous page)

```

Table variable names with their type and shape:
>>> for name in table.colnames:
...     print(name, ' := %s, %s' % (table.coldtypes[name], table.coldtypes[name].
      ↪shape))
ADCCount := uint16, ()
TDCcount := uint8, ()
energy := float64, ()
grid_i := int32, ()
grid_j := int32, ()
idnumber := int64, ()
name := |S16, ()
pressure := float32, ()

```

Here, the name, title, nrow, colnames and coldtypes attributes (see Table for a complete attribute list) of the Table object gives us quite a bit of information about the table data.

You can interactively retrieve general information about the public objects in PyTables by asking for help:

```

>>> help(table)
Help on Table in module tables.table:
class Table(tableextension.Table, tables.leaf.Leaf)
|   This class represents heterogeneous datasets in an HDF5 file.
|
|   Tables are leaves (see the `Leaf` class) whose data consists of a
|   unidimensional sequence of *rows*, where each row contains one or
|   more *fields*. Fields have an associated unique *name* and
|   *position*, with the first field having position 0. All rows have
|   the same fields, which are arranged in *columns*.
[snip]
|
|   Instance variables
|   -----
|
|   The following instance variables are provided in addition to those
|   in `Leaf`. Please note that there are several `col` dictionaries
|   to ease retrieving information about a column directly by its path
|   name, avoiding the need to walk through `Table.description` or
|   `Table.cols`.
|
|   autoindex
|       Automatically keep column indexes up to date?
|
|       Setting this value states whether existing indexes should be
|       automatically updated after an append operation or recomputed
|       after an index-invalidating operation (i.e. removal and
|       modification of rows). The default is true.
[snip]
|   rowsize
|       The size in bytes of each row in the table.
|
|   Public methods -- reading
|   -----
|
|   * col(name)
|   * iterrows([start][, stop][, step])
|   * itersequence(sequence)
|   * itersorted(sortby[, checkCSI][, start][, stop][, step])

```

(continues on next page)

(continued from previous page)

```
| * read([start][, stop][, step][, field][, coords])
| * read_coordinates(coords[, field])
* read_sorted(sortby[, checkCSI][, field][, start][, stop][, step])
| * __getitem__(key)
| * __iter__()
|
| Public methods -- writing
| -----
|
| * append(rows)
| * modify_column([start][, stop][, step][, column][, colname])
[snip]
```

Try getting help with other object docs by yourself:

```
>>> help(h5file)
>>> help(table.remove_rows)
```

To examine metadata in the `/columns/pressure` Array object:

```
>>> pressureObject = h5file.get_node("/columns", "pressure")
>>> print("Info on the object:", repr(pressureObject))
Info on the object: /columns/pressure (Array(3,)) 'Pressure column selection'
  atom := Float64Atom(shape=(), dflt=0.0)
  maindim := 0
  flavor := 'numpy'
  byteorder := 'little'
  chunkshape := None
>>> print("  shape: ==>", pressureObject.shape)
  shape: ==> (3,)
>>> print("  title: ==>", pressureObject.title)
  title: ==> Pressure column selection
>>> print("  atom: ==>", pressureObject.atom)
  atom: ==> Float64Atom(shape=(), dflt=0.0)
```

Observe that we have used the `File.get_node()` method of the `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 one or the other. `File.get_node()` has the advantage that it can get a node from the pathname string (as in this example) and can also act as a filter to show only nodes in a particular location that are instances of class *classname*. In general, however, I consider natural naming to be more elegant and easier to use, especially if you are using the name completion capability present in interactive console. Try this powerful combination of natural naming and completion capabilities present in most Python consoles, and see how pleasant it is to browse the object tree (well, as pleasant as such an activity can be).

If you look at the type attribute of the `pressureObject` object, you can verify that it is a “float64” array. By looking at its shape attribute, you can deduce that the array on disk is unidimensional and has 3 elements. See `Array` or the internal doc strings for the complete `Array` attribute list.

Reading data from Array objects

Once you have found the desired `Array`, use the `read()` method of the `Array` object to retrieve its data:

```
>>> pressureArray = pressureObject.read()
>>> pressureArray
array([ 25.,  36.,  49.])
>>> print("pressureArray is an object of type:", type(pressureArray))
```

(continues on next page)

(continued from previous page)

```

pressureArray is an object of type: <type 'numpy.ndarray'>
>>> nameArray = h5file.root.columns.name.read()
>>> 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

```

You can see that the `Array.read()` method returns an authentic NumPy object for the `pressureObject` instance by looking at the output of the `type()` call. A `read()` of the `nameArray` object instance returns a native Python list (of strings). The type of the object saved is stored as an HDF5 attribute (named `FLAVOR`) for objects on disk. This attribute is then read as Array meta-information (accessible through in 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 later recover them in their original form. See `File.create_array()` for a complete list of supported objects for the Array object class.

1.3.3 Committing data to tables and arrays

We have seen how to create tables and arrays and how to browse both data and metadata in the object tree. Let's examine more closely now one of the most powerful capabilities of PyTables, namely, how to modify already created tables and arrays¹

Appending data to an existing table

Now, let's have a look at how we can add records to an existing table on disk. Let's use our well-known *readout* Table object and 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()

```

It's the same method we used to fill 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 table².

If you look carefully at the code you will see that we have used the `table.row` attribute to create a table row and fill it with the new values. Each time that its `append()` method is called, the actual row is committed to the output buffer and

¹ Appending data to arrays is also supported, but you need to create special objects called *EArray* (see *The EArray class* for more info).

² Note that you can append not only scalar values to tables, but also fully multidimensional array objects.

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 write operation, or else your tables will not be updated!

Let's have a look at some rows in the modified table and verify that our new data has been appended:

```
>>> 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 |
```

Modifying data in tables

Ok, until now, we've been only reading and writing (appending) values to our tables. But there are times that you need to modify your data once you have saved it on disk (this is specially true when you need to modify the real world data to adapt your goals ;). Let's see how we can modify the values that were saved in our existing tables. We will start modifying single cells in the first row of the Particle table:

```
>>> print("Before modif-->", table[0])
Before modif--> (0, 0, 0.0, 0, 10, 0L, 'Particle: 0', 0.0)
>>> table.cols.TDCcount[0] = 1
>>> print("After modifying first row of ADCcount-->", table[0])
After modifying first row of ADCcount--> (0, 1, 0.0, 0, 10, 0L, 'Particle: 0', 0.
↪0)
>>> table.cols.energy[0] = 2
>>> print("After modifying first row of energy-->", table[0])
After modifying first row of energy--> (0, 1, 2.0, 0, 10, 0L, 'Particle: 0', 0.0)
```

We can modify complete ranges of columns as well:

```
>>> table.cols.TDCcount[2:5] = [2,3,4]
>>> print("After modifying slice [2:5] of TDCcount-->", table[0:5])
After modifying slice [2:5] of TDCcount-->
[(0, 1, 2.0, 0, 10, 0L, 'Particle: 0', 0.0)
 (256, 1, 1.0, 1, 9, 17179869184L, 'Particle: 1', 1.0)
 (512, 2, 256.0, 2, 8, 34359738368L, 'Particle: 2', 4.0)
 (768, 3, 6561.0, 3, 7, 51539607552L, 'Particle: 3', 9.0)
 (1024, 4, 65536.0, 4, 6, 68719476736L, 'Particle: 4', 16.0)]
>>> table.cols.energy[1:9:3] = [2,3,4]
>>> print("After modifying slice [1:9:3] of energy-->", table[0:9])
```

(continues on next page)

(continued from previous page)

```

After modifying slice [1:9:3] of energy-->
[(0, 1, 2.0, 0, 10, 0L, 'Particle:      0', 0.0)
 (256, 1, 2.0, 1, 9, 17179869184L, 'Particle:      1', 1.0)
 (512, 2, 256.0, 2, 8, 34359738368L, 'Particle:      2', 4.0)
 (768, 3, 6561.0, 3, 7, 51539607552L, 'Particle:      3', 9.0)
 (1024, 4, 3.0, 4, 6, 68719476736L, 'Particle:      4', 16.0)
 (1280, 5, 390625.0, 5, 5, 85899345920L, 'Particle:      5', 25.0)
 (1536, 6, 1679616.0, 6, 4, 103079215104L, 'Particle:      6', 36.0)
 (1792, 7, 4.0, 7, 3, 120259084288L, 'Particle:      7', 49.0)
 (2048, 8, 16777216.0, 8, 2, 137438953472L, 'Particle:      8', 64.0)]

```

Check that the values have been correctly modified!

Hint: remember that column TDCcount is the second one, and that energy is the third. Look for more info on modifying columns in `Column.__setitem__()`.

PyTables also lets you modify complete sets of rows at the same time. As a demonstration of these capability, see the next example:

```

>>> table.modify_rows(start=1, step=3,
...                    rows=[(1, 2, 3.0, 4, 5, 6L, 'Particle:  None', 8.0),
...                          (2, 4, 6.0, 8, 10, 12L, 'Particle: None*2', 16.0)])
2
>>> print("After modifying the complete third row-->", table[0:5])
After modifying the complete third row-->
[(0, 1, 2.0, 0, 10, 0L, 'Particle:      0', 0.0)
 (1, 2, 3.0, 4, 5, 6L, 'Particle:  None', 8.0)
 (512, 2, 256.0, 2, 8, 34359738368L, 'Particle:      2', 4.0)
 (768, 3, 6561.0, 3, 7, 51539607552L, 'Particle:      3', 9.0)
 (2, 4, 6.0, 8, 10, 12L, 'Particle: None*2', 16.0)]

```

As you can see, the `modify_rows()` call has modified the rows second and fifth, and it returned the number of modified rows.

Apart of `Table.modify_rows()`, there exists another method, called `Table.modify_column()` to modify specific columns as well.

Finally, it exists another way of modifying tables that is generally more handy than the described above. This new way uses the method `Row.update()` of the `Row` instance that is attached to every table, so it is meant to be used in table iterators. Look at the next example:

```

>>> for row in table.where('TDCcount <= 2'):
...     row['energy'] = row['TDCcount']*2
...     row.update()
>>> print("After modifying energy column (where TDCcount <=2)-->", table[0:4])
After modifying energy column (where TDCcount <=2)-->
[(0, 1, 2.0, 0, 10, 0L, 'Particle:      0', 0.0)
 (1, 2, 4.0, 4, 5, 6L, 'Particle:  None', 8.0)
 (512, 2, 4.0, 2, 8, 34359738368L, 'Particle:      2', 4.0)
 (768, 3, 6561.0, 3, 7, 51539607552L, 'Particle:      3', 9.0)]

```

Note: The authors find this way of updating tables (i.e. using `Row.update()`) to be both convenient and efficient. Please make sure to use it extensively.

Caveat emptor: Currently, `Row.update()` will not work (the table will not be updated) if the loop is broken with `break` statement. A possible workaround consists in manually flushing the row internal buffer by calling `row._flushModRows()` just before the `break` statement.

Modifying data in arrays

We are going now to see how to modify data in array objects. The basic way to do this is through the use of `Array.__setitem__()` special method. Let's see at how modify data on the `pressureObject` array:

```
>>> pressureObject = h5file.root.columns.pressure
>>> print("Before modif-->", pressureObject[:])
Before modif--> [ 25.  36.  49.]
>>> pressureObject[0] = 2
>>> print("First modif-->", pressureObject[:])
First modif--> [  2.  36.  49.]
>>> pressureObject[1:3] = [2.1, 3.5]
>>> print("Second modif-->", pressureObject[:])
Second modif--> [ 2.   2.1  3.5]
>>> pressureObject[:,2] = [1,2]
>>> print("Third modif-->", pressureObject[:])
Third modif--> [ 1.   2.1  2. ]
```

So, in general, you can use any combination of (multidimensional) extended slicing.

With the sole exception that you cannot use negative values for step to refer to indexes that you want to modify. See `Array.__getitem__()` for more examples on how to use extended slicing in PyTables objects.

Similarly, with an array of strings:

```
>>> nameObject = h5file.root.columns.name
>>> print("Before modif-->", nameObject[:])
Before modif--> ['Particle:      5', 'Particle:      6', 'Particle:      7']
>>> nameObject[0] = 'Particle:  None'
>>> print("First modif-->", nameObject[:])
First modif--> ['Particle:  None', 'Particle:      6', 'Particle:      7']
>>> nameObject[1:3] = ['Particle:      0', 'Particle:      1']
>>> print("Second modif-->", nameObject[:])
Second modif--> ['Particle:  None', 'Particle:      0', 'Particle:      1']
>>> nameObject[:,2] = ['Particle:     -3', 'Particle:     -5']
>>> print("Third modif-->", nameObject[:])
Third modif--> ['Particle:     -3', 'Particle:      0', 'Particle:     -5']
```

And finally... how to delete rows from a table

We'll finish this tutorial by deleting some rows from the table we have. Suppose that we want to delete the 5th to 9th rows (inclusive):

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

`Table.remove_rows()` 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. Don't forget to close the file when you finish:

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

In [Figure 2](#) you can see a graphical view of the PyTables file with the datasets we have just created. In [Figure 3](#), *General properties of the /detector/readout table*, are displayed the general properties of the table /detector/readout.

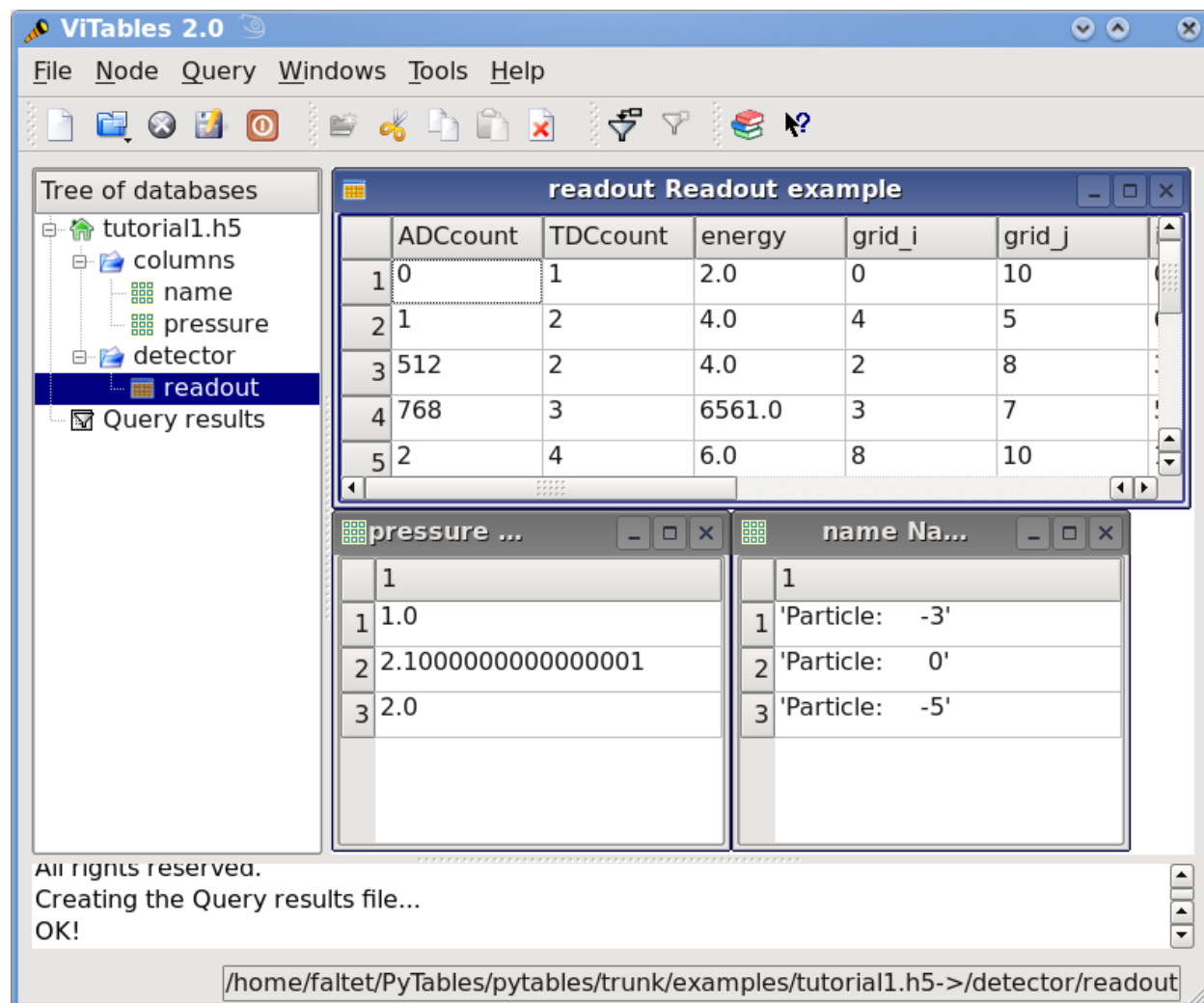


Fig. 4: Figure 2. The final version of the data file for tutorial 1.

1.3.4 Multidimensional table cells and automatic sanity checks

Now it's time for a more real-life example (i.e. with errors in the code). We will create two groups that branch directly from the root node, Particles and Events. Then, we will put three tables in each group. In Particles we will put tables based on the Particle descriptor and in Events, the tables based the Event descriptor.

Afterwards, we will provision the tables with a number of records. Finally, we will read the newly-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 appears to do all of the above, but it contains some small bugs. Note that this Particle class is not directly related to the one defined in last tutorial; this

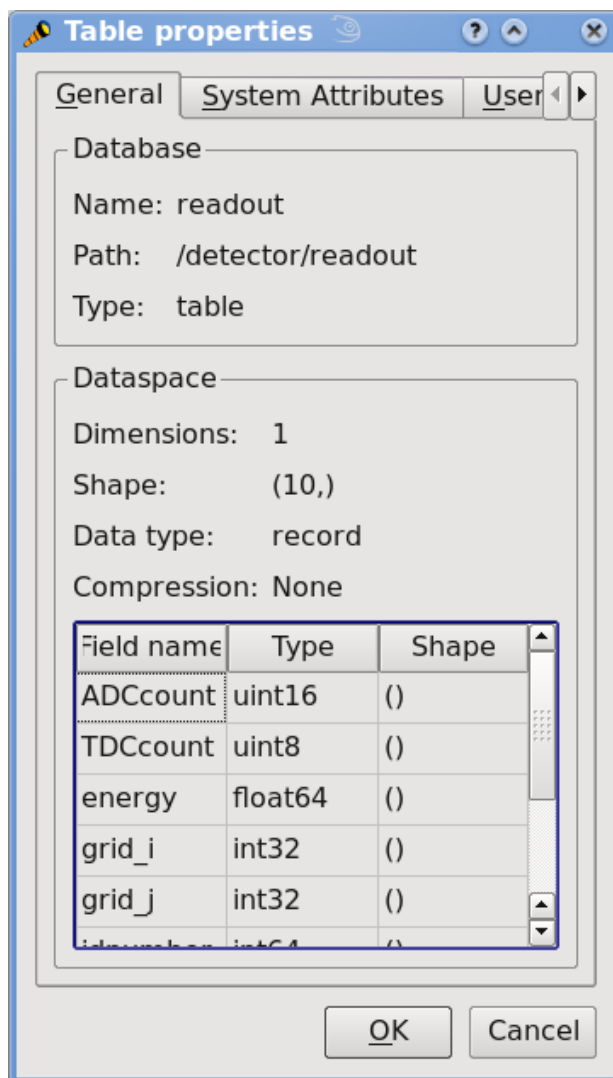


Fig. 5: Figure 3. General properties of the /detector/readout table.

class is simpler (note, however, the *multidimensional* columns called pressure and temperature).

We also introduce a new manner to describe a Table as a structured NumPy dtype (or even as a dictionary), as you can see in the Event description. See `File.create_table()` about the different kinds of descriptor objects that can be passed to this method:

```
from tables import *
from numpy import *

# Describe a particle record
class Particle(IsDescription):
    name          = StringCol(itemsize=16)  # 16-character string
    lati          = Int32Col()              # integer
    longi         = Int32Col()              # integer
    pressure      = Float32Col(shape=(2,3)) # array of floats (single-precision)
    temperature   = Float64Col(shape=(2,3)) # array of doubles (double-precision)

# Native NumPy dtype instances are also accepted
Event = dtype([
    ("name"      , "S16"),
    ("TDCcount"  , uint8),
    ("ADCcount"  , uint16),
    ("xcoord"    , float32),
    ("ycoord"    , float32)
])

# And dictionaries too (this defines the same structure as above)
# Event = {
#     "name"      : StringCol(itemsize=16),
#     "TDCcount"  : UInt8Col(),
#     "ADCcount"  : UInt16Col(),
#     "xcoord"    : Float32Col(),
#     "ycoord"    : Float32Col(),
# }

# Open a file in "w"rite mode
fileh = open_file("tutorial2.h5", mode = "w")

# Get the HDF5 root group
root = fileh.root

# Create the groups:
for groupname in ("Particles", "Events"):
    group = fileh.create_group(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.create_table("/Particles", tablename, Particle, "Particles:
↪"+tablename)

    # Get the record object associated with the table:
    particle = table.row

    # Fill the table with 257 particles
```

(continues on next page)

(continued from previous page)

```

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*xrange(2*3)).reshape((2,4)) # Incorrect
    #particle['pressure'] = array(i*xrange(2*3)).reshape((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.create_table(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['xcoor'] = float(i**2)     # Wrong spelling
        #event['xcoord'] = float(i**2)   # Correct spelling
        event['ADCCount'] = "sss"        # Wrong type
        #event['ADCCount'] = i * 2       # Correct 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()

```

Shape checking

If you look at the code carefully, you'll see that it won't work. You will get the following error.

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 60, in <module>
    particle['pressure'] = array(i*arange(2*3)).reshape((2,4)) # Incorrect
ValueError: total size of new array must be unchanged
Closing remaining open files: tutorial2.h5... done
```

This error indicates that you are trying to assign an array with an incompatible shape to a table cell. Looking at the source, we see that we were trying to assign an array of shape (2,4) to a pressure element, which was defined with the shape (2,3).

In general, these kinds of operations are forbidden, with one valid exception: when you assign a *scalar* value to a multidimensional column cell, all the cell elements are populated with the value of the scalar. For example:

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

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

Field name checking

After fixing the previous error and rerunning the program, we encounter another error.

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 73, in ?
    event['xcoor'] = float(i**2) # Wrong spelling
  File "tableextension.pyx", line 1094, in tableextension.Row.__setitem__
  File "tableextension.pyx", line 127, in tableextension.get_nested_field_cache
  File "utilsextension.pyx", line 331, in utilsextension.get_nested_field
KeyError: 'no such column: xcoor'
```

This error indicates that we are attempting 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 unusual behavior for Python, as normally when you assign a value to a non-existent instance variable, Python creates a new variable with that name. Such a feature can be dangerous when dealing with an object that contains a fixed list of field names. PyTables checks that the field exists and raises a *KeyError* if the check fails.

Data type checking

Finally, the last issue which we will find here is a *TypeError* exception.

```
$ python tutorial2.py
Traceback (most recent call last):
  File "tutorial2.py", line 75, in ?
    event['ADCCount'] = "sss" # Wrong type
  File "tableextension.pyx", line 1111, in tableextension.Row.__setitem__
TypeError: invalid type (<type 'str'>) for column ``ADCCount``
```

And, if we change the affected line to read:

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

we will see that the script ends well.

You can see the structure created with this (corrected) script in [Figure 4](#). In particular, note the multidimensional column cells in table /Particles/TParticle2.

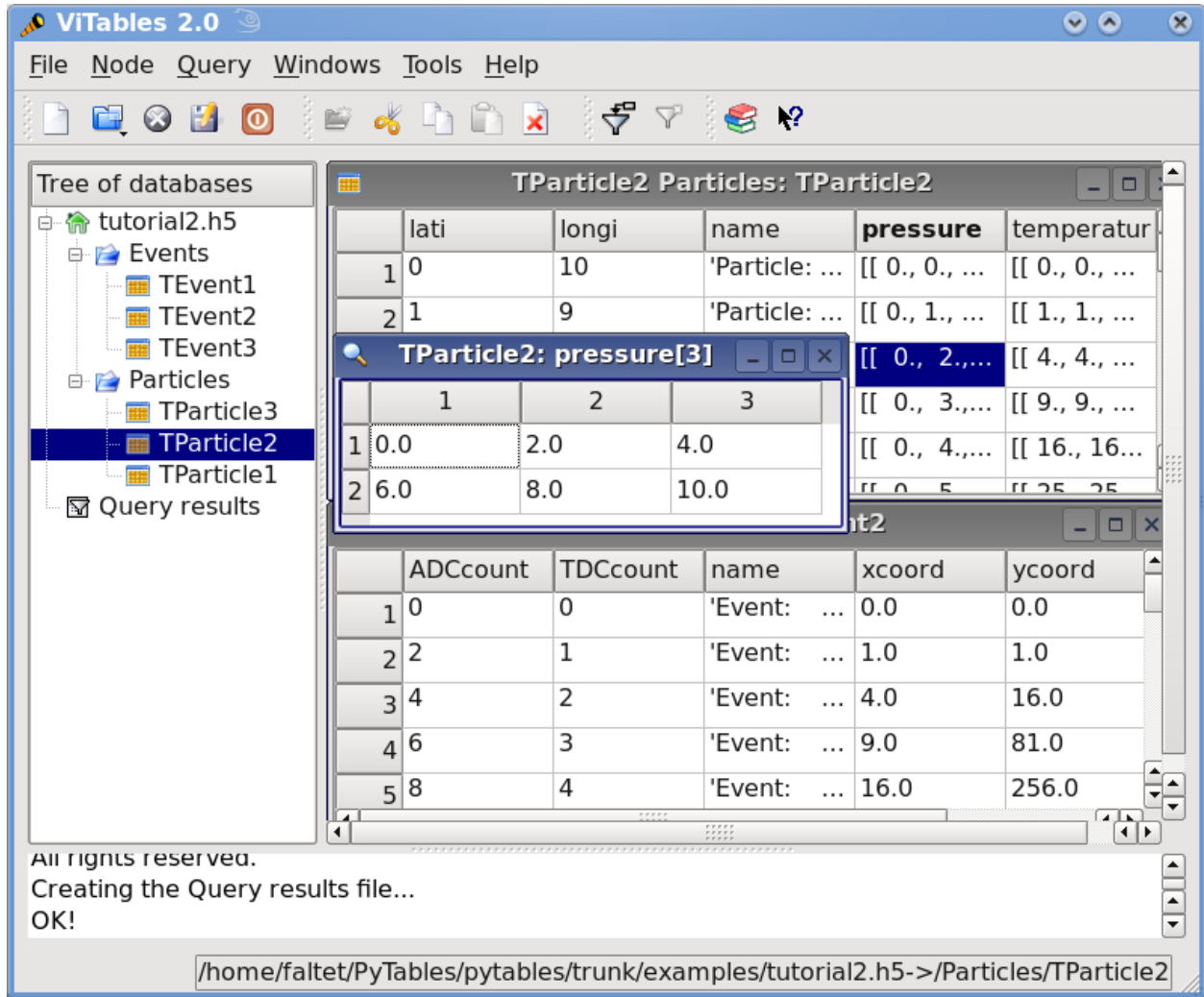


Fig. 6: **Figure 4.** Table hierarchy for tutorial 2.

1.3.5 Using links for more convenient access to nodes

Links are special nodes that can be used to create additional paths to your existing nodes. PyTables supports three kinds of links: hard links, soft links (aka symbolic links) and external links.

Hard links let the user create additional paths to access another node in the same file, and once created, they are indistinguishable from the referred node object, except that they have different paths in the object tree. For example, if the referred node is, say, a Table object, then the new hard link will become a Table object itself. From this point on, you will be able to access the same Table object from two different paths: the original one and the new hard link path. If you delete one path to the table, you will be able to reach it via the other path.

Soft links are similar to hard links, but they keep their own personality. When you create a soft link to another node, you will get a new `SoftLink` object that *refers* to that node. However, in order to access the referred node, you need to *dereference* it.

Finally, external links are like soft links, with the difference that these are meant to point to nodes in *external* files instead of nodes in the same file. They are represented by the `ExternalLink` class and, like soft links, you need to dereference them in order to get access to the pointed node.

Interactive example

Now we are going to learn how to deal with links. You can find the code used in this section in `examples/links.py`.

First, let's create a file with some group structure:

```
>>> import tables as tb
>>> f1 = tb.open_file('links1.h5', 'w')
>>> g1 = f1.create_group('/', 'g1')
>>> g2 = f1.create_group(g1, 'g2')
```

Now, we will put some datasets on the `/g1` and `/g1/g2` groups:

```
>>> a1 = f1.create_carray(g1, 'a1', tb.Int64Atom(), shape=(10000,))
>>> t1 = f1.create_table(g2, 't1', {'f1': tb.IntCol(), 'f2': tb.FloatCol()})
```

We can start the party now. We are going to create a new group, say `/gl`, where we will put our links and will start creating one hard link too:

```
>>> gl = f1.create_group('/', 'gl')
>>> ht = f1.create_hard_link(gl, 'ht', '/g1/g2/t1') # ht points to t1
>>> print("`%s` is a hard link to: `%s`" % (ht, t1))
`/gl/ht (Table(0,))` is a hard link to: `/g1/g2/t1 (Table(0,))`
```

You can see how we've created a hard link in `/gl/ht` which is pointing to the existing table in `/g1/g2/t1`. Have look at how the hard link is represented; it looks like a table, and actually, it is an *real* table. We have two different paths to access that table, the original `/g1/g2/t1` and the new one `/gl/ht`. If we remove the original path we still can reach the table by using the new path:

```
>>> t1.remove()
>>> print("table continues to be accessible in: `%s`" % f1.get_node('/gl/ht'))
table continues to be accessible in: `/gl/ht (Table(0,))`
```

So far so good. Now, let's create a couple of soft links:

```
>>> la1 = f1.create_soft_link(gl, 'la1', '/g1/a1') # la1 points to a1
>>> print("`%s` is a soft link to: `%s`" % (la1, la1.target))
`/gl/la1 (SoftLink) -> /g1/a1` is a soft link to: `/g1/a1`
>>> lt = f1.create_soft_link(gl, 'lt', '/g1/g2/t1') # lt points to t1
>>> print("`%s` is a soft link to: `%s`" % (lt, lt.target))
`/gl/lt (SoftLink) -> /g1/g2/t1 (dangling)` is a soft link to: `/g1/g2/t1`
```

Okay, we see how the first link `/gl/la1` points to the array `/g1/a1`. Notice how the link prints as a `SoftLink`, and how the referred node is stored in the `target` instance attribute. The second link (`/gl/lt`) pointing to `/g1/g2/t1` also has been created successfully, but by better inspecting the string representation of it, we see that is labeled as '(dangling)'. Why is this? Well, you should remember that we recently removed the `/g1/g2/t1` path to access table `t1`. When printing it, the object knows that it points to *nowhere* and reports this. This is a nice way to quickly know whether a soft link points to an existing node or not.

So, let's re-create the removed path to t1 table:

```
>>> t1 = f1.create_hard_link('/g1/g2', 't1', '/g1/ht')
>>> print("`%s` is not dangling anymore" % (lt,))
`/g1/lt (SoftLink) -> /g1/g2/t1` is not dangling anymore
```

and the soft link is pointing to an existing node now.

Of course, for soft links to serve any actual purpose we need a way to get the pointed node. It happens that soft links are callable, and that's the way to get the referred nodes back:

```
>>> plt = lt()
>>> print("dereferenced lt node: `%s`" % plt)
dereferenced lt node: `/g1/g2/t1 (Table(0,)) `
>>> pla1 = la1()
>>> print("dereferenced la1 node: `%s`" % pla1)
dereferenced la1 node: `/g1/a1 (CArray(10000,)) `
```

Now, plt is a Python reference to the t1 table while pla1 refers to the a1 array. Easy, uh?

Let's suppose now that a1 is an array whose access speed is critical for our application. One possible solution is to move the entire file into a faster disk, say, a solid state disk so that access latencies can be reduced quite a lot. However, it happens that our file is too big to fit into our shiny new (although small in capacity) SSD disk. A solution is to copy just the a1 array into a separate file that would fit into our SSD disk. However, our application would be able to handle two files instead of only one, adding significantly more complexity, which is not a good thing.

External links to the rescue! As we've already said, external links are like soft links, but they are designed to link objects in external files. Back to our problem, let's copy the a1 array into a different file:

```
>>> f2 = tb.open_file('links2.h5', 'w')
>>> new_a1 = a1.copy(f2.root, 'a1')
>>> f2.close() # close the other file
```

And now, we can remove the existing soft link and create the external link in its place:

```
>>> la1.remove()
>>> la1 = f1.create_external_link(g1, 'la1', 'links2.h5:/a1')
>>> print("`%s` is an external link to: `%s`" % (la1, la1.target))
`/g1/la1 (ExternalLink) -> links2.h5:/a1` is an external link to: `links2.h5:/a1`
```

Let's try dereferencing it:

```
>>> new_a1 = la1() # dereferencing la1 returns a1 in links2.h5
>>> print("dereferenced la1 node: `%s`" % new_a1)
dereferenced la1 node: `/a1 (CArray(10000,)) `
```

Well, it seems like we can access the external node. But just to make sure that the node is in the other file:

```
>>> print("new_a1 file:", new_a1._v_file.filename)
new_a1 file: links2.h5
```

Okay, the node is definitely in the external file. So, you won't have to worry about your application: it will work exactly the same no matter the link is internal (soft) or external.

Finally, here it is a dump of the objects in the final file, just to get a better idea of what we ended with:

```
>>> f1.close()
>>> exit()
```

(continues on next page)

(continued from previous page)

```
$ ptdump links1.h5
/ (RootGroup) ''
/g1 (Group) ''
/g1/a1 (CArray(10000,)) ''
/g1 (Group) ''
/g1/ht (Table(0,)) ''
/g1/la1 (ExternalLink) -> links2.h5:/a1
/g1/lt (SoftLink) -> /g1/g2/t1
/g1/g2 (Group) ''
/g1/g2/t1 (Table(0,)) ''
```

This ends this tutorial. I hope it helped you to appreciate how useful links can be. I'm sure you will find other ways in which you can use links that better fit your own needs.

1.3.6 Exercising the Undo/Redo feature

PyTables has integrated support for undoing and/or redoing actions. This functionality lets you put marks in specific places of your hierarchy manipulation operations, so that you can make your HDF5 file pop back (*undo*) to a specific mark (for example for inspecting how your hierarchy looked at that point). You can also go forward to a more recent marker (*redo*). You can even do jumps to the marker you want using just one instruction as we will see shortly.

You can undo/redo all the operations that are related to object tree management, like creating, deleting, moving or renaming nodes (or complete sub-hierarchies) inside a given object tree. You can also undo/redo operations (i.e. creation, deletion or modification) of persistent node attributes. However, when actions include *internal* modifications of datasets (that includes `Table.append`, `Table.modify_rows` or `Table.remove_rows` among others), they cannot be undone/redone currently.

This capability can be useful in many situations, like for example when doing simulations with multiple branches. When you have to choose a path to follow in such a situation, you can put a mark there and, if the simulation is not going well, you can go back to that mark and start another path. Other possible application is defining coarse-grained operations which operate in a transactional-like way, i.e. which return the database to its previous state if the operation finds some kind of problem while running. You can probably devise many other scenarios where the Undo/Redo feature can be useful to you³.

A basic example

In this section, we are going to show the basic behavior of the Undo/Redo feature. You can find the code used in this example in `examples/tutorial3-1.py`. A somewhat more complex example will be explained in the next section.

First, let's create a file:

```
>>> import tables
>>> fileh = tables.open_file("tutorial3-1.h5", "w", title="Undo/Redo demo 1")
```

And now, activate the Undo/Redo feature with the method `File.enable_undo()` of `File`:

```
>>> fileh.enable_undo()
```

From now on, all our actions will be logged internally by PyTables. Now, we are going to create a node (in this case an Array object):

³ You can even *hide* nodes temporarily. Will you be able to find out how?

```
>>> one = fileh.create_array('/', 'anarray', [3,4], "An array")
```

Now, mark this point:

```
>>> fileh.mark()
1
```

We have marked the current point in the sequence of actions. In addition, the `mark()` method has returned the identifier assigned to this new mark, that is 1 (mark #0 is reserved for the implicit mark at the beginning of the action log). In the next section we will see that you can also assign a *name* to a mark (see `File.mark()` for more info on `mark()`). Now, we are going to create another array:

```
>>> another = fileh.create_array('/', 'anotherarray', [4,5], "Another array")
```

Right. Now, we can start doing funny things. Let's say that we want to pop back to the previous mark (that whose value was 1, do you remember?). Let's introduce the `undo()` method (see `File.undo()`):

```
>>> fileh.undo()
```

Fine, what do you think it happened? Well, let's have a look at the object tree:

```
>>> print(fileh)
tutorial3-1.h5 (File) 'Undo/Redo demo 1'
Last modif.: 'Tue Mar 13 11:43:55 2007'
Object Tree:
/ (RootGroup) 'Undo/Redo demo 1'
/anarray (Array(2,)) 'An array'
```

What happened with the `/anotherarray` node we've just created? You guess it, it has disappeared because it was created *after* the mark 1. If you are curious enough you may well ask where it has gone. Well, it has not been deleted completely; it has been just moved into a special, hidden, group of PyTables that renders it invisible and waiting for a chance to be reborn.

Now, unwind once more, and look at the object tree:

```
>>> fileh.undo()
>>> print(fileh)
tutorial3-1.h5 (File) 'Undo/Redo demo 1'
Last modif.: 'Tue Mar 13 11:43:55 2007'
Object Tree:
/ (RootGroup) 'Undo/Redo demo 1'
```

Oops, `/anarray` has disappeared as well!. Don't worry, it will revisit us very shortly. So, you might be somewhat lost right now; in which mark are we?. Let's ask the `File.get_current_mark()` method in the file handler:

```
>>> print(fileh.get_current_mark())
0
```

So we are at mark #0, remember? Mark #0 is an implicit mark that is created when you start the log of actions when calling `File.enable_undo()`. Fine, but you are missing your too-young-to-die arrays. What can we do about that? `File.redo()` to the rescue:

```
>>> fileh.redo()
>>> print(fileh)
tutorial3-1.h5 (File) 'Undo/Redo demo 1'
Last modif.: 'Tue Mar 13 11:43:55 2007'
```

(continues on next page)

(continued from previous page)

```
Object Tree:
/ (RootGroup) 'Undo/Redo demo 1'
/anarray (Array(2,)) 'An array'
```

Great! The /anarray array has come into life again. Just check that it is alive and well:

```
>>> fileh.root.anarray.read()
[3, 4]
>>> fileh.root.anarray.title
'An array'
```

Well, it looks pretty similar than in its previous life; what's more, it is exactly the same object!:

```
>>> fileh.root.anarray is one
True
```

It just was moved to the the hidden group and back again, but that's all! That's kind of fun, so we are going to do the same with /anotherarray:

```
>>> fileh.redo()
>>> print(fileh)
tutorial3-1.h5 (File) 'Undo/Redo demo 1'
Last modif.: 'Tue Mar 13 11:43:55 2007'
Object Tree:
/ (RootGroup) 'Undo/Redo demo 1'
/anarray (Array(2,)) 'An array'
/anotherarray (Array(2,)) 'Another array'
```

Welcome back, /anotherarray! Just a couple of sanity checks:

```
>>> assert fileh.root.anotherarray.read() == [4,5]
>>> assert fileh.root.anotherarray.title == "Another array"
>>> fileh.root.anotherarray is another
True
```

Nice, you managed to turn your data back into life. Congratulations! But wait, do not forget to close your action log when you don't need this feature anymore:

```
>>> fileh.disable_undo()
```

That will allow you to continue working with your data without actually requiring PyTables to keep track of all your actions, and more importantly, allowing your objects to die completely if they have to, not requiring to keep them anywhere, and hence saving process time and space in your database file.

A more complete example

Now, time for a somewhat more sophisticated demonstration of the Undo/Redo feature. In it, several marks will be set in different parts of the code flow and we will see how to jump between these marks with just one method call. You can find the code used in this example in `examples/tutorial3-2.py`

Let's introduce the first part of the code:

```
import tables

# Create an HDF5 file
```

(continues on next page)

(continued from previous page)

```

fileh = tables.open_file('tutorial3-2.h5', 'w', title='Undo/Redo demo 2')

    #'-***-***-***-***-***- enable undo/redo log  -***-***-***-***-***-'
fileh.enable_undo()

# Start undoable operations
fileh.create_array('/', 'otherarray1', [3,4], 'Another array 1')
fileh.create_group('/', 'agroup', 'Group 1')

# Create a 'first' mark
fileh.mark('first')
fileh.create_array('/agroup', 'otherarray2', [4,5], 'Another array 2')
fileh.create_group('/agroup', 'agroup2', 'Group 2')

# Create a 'second' mark
fileh.mark('second')
fileh.create_array('/agroup/agroup2', 'otherarray3', [5,6], 'Another array 3')

# Create a 'third' mark
fileh.mark('third')
fileh.create_array('/', 'otherarray4', [6,7], 'Another array 4')
fileh.create_array('/agroup', 'otherarray5', [7,8], 'Another array 5')

```

You can see how we have set several marks interspersed in the code flow, representing different states of the database. Also, note that we have assigned *names* to these marks, namely ‘first’, ‘second’ and ‘third’.

Now, start doing some jumps back and forth in the states of the database:

```

# Now go to mark 'first'
fileh.goto('first')
assert '/otherarray1' in fileh
assert '/agroup' in fileh
assert '/agroup/agroup2' not in fileh
assert '/agroup/otherarray2' not in fileh
assert '/agroup/agroup2/otherarray3' not in fileh
assert '/otherarray4' not in fileh
assert '/agroup/otherarray5' not in fileh

# Go to mark 'third'
fileh.goto('third')
assert '/otherarray1' in fileh
assert '/agroup' in fileh
assert '/agroup/agroup2' in fileh
assert '/agroup/otherarray2' in fileh
assert '/agroup/agroup2/otherarray3' in fileh
assert '/otherarray4' not in fileh
assert '/agroup/otherarray5' not in fileh

# Now go to mark 'second'
fileh.goto('second')
assert '/otherarray1' in fileh
assert '/agroup' in fileh
assert '/agroup/agroup2' in fileh
assert '/agroup/otherarray2' in fileh
assert '/agroup/agroup2/otherarray3' not in fileh
assert '/otherarray4' not in fileh
assert '/agroup/otherarray5' not in fileh

```

Well, the code above shows how easy is to jump to a certain mark in the database by using the `File.goto()` method.

There are also a couple of implicit marks for going to the beginning or the end of the saved states: 0 and -1. Going to mark #0 means go to the beginning of the saved actions, that is, when method `fileh.enable_undo()` was called. Going to mark #-1 means go to the last recorded action, that is the last action in the code flow.

Let's see what happens when going to the end of the action log:

```
# Go to the end
fileh.goto(-1)
assert '/otherarray1' in fileh
assert '/agroup' in fileh
assert '/agroup/agroup2' in fileh
assert '/agroup/otherarray2' in fileh
assert '/agroup/agroup2/otherarray3' in fileh
assert '/otherarray4' in fileh
assert '/agroup/otherarray5' in fileh

# Check that objects have come back to life in a sane state
assert fileh.root.otherarray1.read() == [3,4]
assert fileh.root.agroup.otherarray2.read() == [4,5]
assert fileh.root.agroup.agroup2.otherarray3.read() == [5,6]
assert fileh.root.otherarray4.read() == [6,7]
assert fileh.root.agroup.otherarray5.read() == [7,8]
```

Try yourself going to the beginning of the action log (remember, the mark #0) and check the contents of the object tree.

We have nearly finished this demonstration. As always, do not forget to close the action log as well as the database:

```
#'-*-*-*-*-*-*-*-*-*- disable undo/redo log  -*-*-*-*-*-*-*-*-*-*-
fileh.disable_undo()
# Close the file
fileh.close()
```

You might want to check other examples on Undo/Redo feature that appear in `examples/undo-redo.py`.

1.3.7 Using enumerated types

PyTables includes support for handling enumerated types. Those types are defined by providing an exhaustive *set* or *list* of possible, named values for a variable of that type. Enumerated variables of the same type are usually compared between them for equality and sometimes for order, but are not usually operated upon.

Enumerated values have an associated *name* and *concrete value*. Every name is unique and so are concrete values. An enumerated variable always takes the concrete value, not its name. Usually, the concrete value is not used directly, and frequently it is entirely irrelevant. For the same reason, an enumerated variable is not usually compared with concrete values out of its enumerated type. For that kind of use, standard variables and constants are more adequate.

PyTables provides the Enum (see *The Enum class*) class to provide support for enumerated types. Each instance of Enum is an enumerated type (or *enumeration*). For example, let us create an enumeration of colors

All these examples can be found in `examples/play-with-enums.py`:

```
>>> import tables
>>> colorList = ['red', 'green', 'blue', 'white', 'black']
>>> colors = tables.Enum(colorList)
```

Here we used a simple list giving the names of enumerated values, but we left the choice of concrete values up to the Enum class. Let us see the enumerated pairs to check those values:

```
>>> print("Colors:", [v for v in colors])
Colors: [('blue', 2), ('black', 4), ('white', 3), ('green', 1), ('red', 0)]
```

Names have been given automatic integer concrete values. We can iterate over the values in an enumeration, but we will usually be more interested in accessing single values. We can get the concrete value associated with a name by accessing it as an attribute or as an item (the later can be useful for names not resembling Python identifiers):

```
>>> print("Value of 'red' and 'white':", (colors.red, colors.white))
Value of 'red' and 'white': (0, 3)
>>> print("Value of 'yellow':", colors.yellow)
Value of 'yellow':
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File ".../tables/misc/enum.py", line 230, in __getattr__
    raise AttributeError(*ke.args)
AttributeError: no enumerated value with that name: 'yellow'
>>>
>>> print("Value of 'red' and 'white':", (colors['red'], colors['white']))
Value of 'red' and 'white': (0, 3)
>>> print("Value of 'yellow':", colors['yellow'])
Value of 'yellow':
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File ".../tables/misc/enum.py", line 189, in __getitem__
    raise KeyError("no enumerated value with that name: %r" % (name,))
KeyError: "no enumerated value with that name: 'yellow'"
```

See how accessing a value that is not in the enumeration raises the appropriate exception. We can also do the opposite action and get the name that matches a concrete value by using the `__call__()` method of Enum:

```
>>> print("Name of value %s:" % colors.red, colors(colors.red))
Name of value 0: red
>>> print("Name of value 1234:", colors(1234))
Name of value 1234:
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
  File ".../tables/misc/enum.py", line 320, in __call__
    raise ValueError(
ValueError: no enumerated value with that concrete value: 1234
```

You can see what we made as using the enumerated type to *convert* a concrete value into a name in the enumeration. Of course, values out of the enumeration can not be converted.

Enumerated columns

Columns of an enumerated type can be declared by using the `EnumCol` (see *The Col class and its descendants*) class. To see how this works, let us open a new PyTables file and create a table to collect the simulated results of a probabilistic experiment. In it, we have a bag full of colored balls; we take a ball out and annotate the time of extraction and the color of the ball:

```
>>> h5f = tables.open_file('enum.h5', 'w')
>>> class BallExt(tables.IsDescription):
...     ballTime = tables.Time32Col()
...     ballColor = tables.EnumCol(colors, 'black', base='uint8')
```

(continues on next page)

(continued from previous page)

```
>>> tbl = h5f.create_table('/', 'extractions', BallExt, title="Random ball extractions
↳")
>>>
```

We declared the ballColor column to be of the enumerated type colors, with a default value of black. We also stated that we are going to store concrete values as unsigned 8-bit integer values⁴.

Let us use some random values to fill the table:

```
>>> import time
>>> import random
>>> now = time.time()
>>> row = tbl.row
>>> for i in range(10):
...     row['ballTime'] = now + i
...     row['ballColor'] = colors[random.choice(colorList)] # notice this
...     row.append()
>>>
```

Notice how we used the `__getitem__()` call of colors to get the concrete value to store in ballColor. You should know that this way of appending values to a table does automatically check for the validity on enumerated values. For instance:

```
>>> row['ballTime'] = now + 42
>>> row['ballColor'] = 1234
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "tableextension.pyx", line 1086, in tableextension.Row.__setitem__
  File ".../tables/misc/enum.py", line 320, in __call__
    "no enumerated value with that concrete value: %r" % (value,)
ValueError: no enumerated value with that concrete value: 1234
```

But take care that this check is *only* performed here and not in other methods such as `tbl.append()` or `tbl.modify_rows()`. Now, after flushing the table we can see the results of the insertions:

```
>>> tbl.flush()
>>> for r in tbl:
...     ballTime = r['ballTime']
...     ballColor = colors(r['ballColor']) # notice this
...     print("Ball extracted on %d is of color %s." % (ballTime, ballColor))
Ball extracted on 1173785568 is of color green.
Ball extracted on 1173785569 is of color black.
Ball extracted on 1173785570 is of color white.
Ball extracted on 1173785571 is of color black.
Ball extracted on 1173785572 is of color black.
Ball extracted on 1173785573 is of color red.
Ball extracted on 1173785574 is of color green.
Ball extracted on 1173785575 is of color red.
Ball extracted on 1173785576 is of color white.
Ball extracted on 1173785577 is of color white.
```

As a last note, you may be wondering how to have access to the enumeration associated with ballColor once the file is closed and reopened. You can call `tbl.get_enum('ballColor')` (see `Table.get_enum()`) to get the enumeration back.

⁴ In fact, only integer values are supported right now, but this may change in the future.

Enumerated arrays

EArray and VArray leaves can also be declared to store enumerated values by means of the EnumAtom (see *The Atom class and its descendants*) class, which works very much like EnumCol for tables. Also, Array leaves can be used to open native HDF enumerated arrays.

Let us create a sample EArray containing ranges of working days as bidimensional values:

```
>>> workingDays = {'Mon': 1, 'Tue': 2, 'Wed': 3, 'Thu': 4, 'Fri': 5}
>>> dayRange = tables.EnumAtom(workingDays, 'Mon', base='uint16')
>>> earr = h5f.create_earray('/', 'days', dayRange, (0, 2), title="Working day ranges
↪")
>>> earr.flavor = 'python'
```

Nothing surprising, except for a pair of details. In the first place, we use a *dictionary* instead of a list to explicitly set concrete values in the enumeration. In the second place, there is no explicit Enum instance created! Instead, the dictionary is passed as the first argument to the constructor of EnumAtom. If the constructor gets a list or a dictionary instead of an enumeration, it automatically builds the enumeration from it.

Now let us feed some data to the array:

```
>>> wdays = earr.get_enum()
>>> earr.append([(wdays.Mon, wdays.Fri), (wdays.Wed, wdays.Fri)])
>>> earr.append([(wdays.Mon, 1234)])
```

Please note that, since we had no explicit Enum instance, we were forced to use get_enum() (see *EArray methods*) to get it from the array (we could also have used dayRange.enum). Also note that we were able to append an invalid value (1234). Array methods do not check the validity of enumerated values.

Finally, we will print the contents of the array:

```
>>> for (d1, d2) in earr:
...     print("From %s to %s (%d days)." % (wdays(d1), wdays(d2), d2-d1+1))
From Mon to Fri (5 days).
From Wed to Fri (3 days).
Traceback (most recent call last):
  File "<stdin>", line 2, in <module>
  File ".../tables/misc/enum.py", line 320, in __call__
    "no enumerated value with that concrete value: %r" % (value,))
ValueError: no enumerated value with that concrete value: 1234
```

That was an example of operating on concrete values. It also showed how the value-to-name conversion failed because of the value not belonging to the enumeration.

Now we will close the file, and this little tutorial on enumerated types is done:

```
>>> h5f.close()
```

1.3.8 Dealing with nested structures in tables

PyTables supports the handling of nested structures (or nested datatypes, as you prefer) in table objects, allowing you to define arbitrarily nested columns.

An example will clarify what this means. Let's suppose that you want to group your data in pieces of information that are more related than others pieces in your table. So you may want to tie them up together in order to have your table better structured but also be able to retrieve and deal with these groups more easily.

You can create such a nested substructures by just nesting subclasses of `IsDescription`. Let's see one example (okay, it's a bit silly, but will serve for demonstration purposes):

```
from tables import *

class Info(IsDescription):
    """A sub-structure of Test"""
    _v_pos = 2 # The position in the whole structure
    name = StringCol(10)
    value = Float64Col(pos=0)

colors = Enum(['red', 'green', 'blue'])

class NestedDescr(IsDescription):
    """A description that has several nested columns"""
    color = EnumCol(colors, 'red', base='uint32')
    info1 = Info()

    class info2(IsDescription):
        _v_pos = 1
        name = StringCol(10)
        value = Float64Col(pos=0)

        class info3(IsDescription):
            x = Float64Col(dflt=1)
            y = UInt8Col(dflt=1)
```

The root class is `NestedDescr` and both `info1` and `info2` are *substructures* of it. Note how `info1` is actually an instance of the class `Info` that was defined prior to `NestedDescr`. Also, there is a third substructure, namely `info3` that hangs from the substructure `info2`. You can also define positions of substructures in the containing object by declaring the special class attribute `_v_pos`.

Nested table creation

Now that we have defined our nested structure, let's create a *nested* table, that is a table with columns that contain other subcolumns:

```
>>> fileh = open_file("nested-tut.h5", "w")
>>> table = fileh.create_table(fileh.root, 'table', NestedDescr)
```

Done! Now, we have to feed the table with some values. The problem is how we are going to reference to the nested fields. That's easy, just use a `'/'` character to separate names in different nested levels. Look at this:

```
>>> row = table.row
>>> for i in range(10):
...     row['color'] = colors[['red', 'green', 'blue'][i%3]]
...     row['info1/name'] = "name1-%s" % i
...     row['info2/name'] = "name2-%s" % i
...     row['info2/info3/y'] = i
...     # All the rest will be filled with defaults
...     row.append()
>>> table.flush()
>>> table.nrows
10
```

You see? In order to fill the fields located in the substructures, we just need to specify its full path in the table hierarchy.

Reading nested tables

Now, what happens if we want to read the table? What kind of data container will we get? Well, it's worth trying it:

```
>>> nra = table[:,4]
>>> nra
array([((1.0, 0), 'name2-0', 0.0), ('name1-0', 0.0), 0L),
      ((1.0, 4), 'name2-4', 0.0), ('name1-4', 0.0), 1L),
      ((1.0, 8), 'name2-8', 0.0), ('name1-8', 0.0), 2L)],
      dtype=[('info2', [('info3', [('x', '>f8'), ('y', '\\|u1')]),
        ('name', '\\|S10'), ('value', '>f8')]),
        ('info1', [('name', '\\|S10'), ('value', '>f8')]),
        ('color', '>u4')])
```

What we got is a NumPy array with a *compound, nested datatype* (its dtype is a list of name-datatype tuples). We read one row for each four in the table, giving a result of three rows.

You can make use of the above object in many different ways. For example, you can use it to append new data to the existing table object:

```
>>> table.append(nra)
>>> table.nrows
13
```

Or, to create new tables:

```
>>> table2 = fileh.create_table(fileh.root, 'table2', nra)
>>> table2[:]
array([((1.0, 0), 'name2-0', 0.0), ('name1-0', 0.0), 0L),
      ((1.0, 4), 'name2-4', 0.0), ('name1-4', 0.0), 1L),
      ((1.0, 8), 'name2-8', 0.0), ('name1-8', 0.0), 2L)],
      dtype=[('info2', [('info3', [('x', '<f8'), ('y', '\\|u1')]),
        ('name', '\\|S10'), ('value', '<f8')]),
        ('info1', [('name', '\\|S10'), ('value', '<f8')]),
        ('color', '<u4')])
```

Finally, we can select nested values that fulfill some condition:

```
>>> names = [ x['info2/name'] for x in table if x['color'] == colors.red ]
>>> names
['name2-0', 'name2-3', 'name2-6', 'name2-9', 'name2-0']
```

Note that the row accessor does not provide the natural naming feature, so you have to completely specify the path of your desired columns in order to reach them.

Using Cols accessor

We can use the cols attribute object (see *The Cols class*) of the table so as to quickly access the info located in the interesting substructures:

```
>>> table.cols.info2[1:5]
array([(1.0, 1), 'name2-1', 0.0), ((1.0, 2), 'name2-2', 0.0),
      ((1.0, 3), 'name2-3', 0.0), ((1.0, 4), 'name2-4', 0.0)],
      dtype=[('info3', [('x', '<f8'), ('y', '\\|u1')]), ('name', '\\|S10'),
        ('value', '<f8')])
```

Here, we have made use of the `cols` accessor to access to the *info2* substructure and an slice operation to get access to the subset of data we were interested in; you probably have recognized the natural naming approach here. We can continue and ask for data in *info3* substructure:

```
>>> table.cols.info2.info3[1:5]
array([(1.0, 1), (1.0, 2), (1.0, 3), (1.0, 4)],
      dtype=[('x', '<f8'), ('y', '\u1')])
```

You can also use the `_f_col` method to get a handler for a column:

```
>>> table.cols._f_col('info2')
/table.cols.info2 (Cols), 3 columns
  info3 (Cols(), Description)
    name (Column(), \S10)
    value (Column(), float64)
```

Here, you've got another `Cols` object handler because *info2* was a nested column. If you select a non-nested column, you will get a regular `Column` instance:

```
>>> table.cols._f_col('info2/info3/y')
/table.cols.info2.info3.y (Column(), uint8, idx=None)
```

To sum up, the `cols` accessor is a very handy and powerful way to access data in your nested tables. Don't be afraid of using it, specially when doing interactive work.

Accessing meta-information of nested tables

Tables have an attribute called `description` which points to an instance of the *Description* class (see *The Description class*) and is useful to discover different meta-information about table data.

Let's see how it looks like:

```
>>> table.description
{
  "info2": {
    "info3": {
      "x": Float64Col(shape=(), dflt=1.0, pos=0),
      "y": UInt8Col(shape=(), dflt=1, pos=1)},
    "name": StringCol(itemsize=10, shape=(), dflt='', pos=1),
    "value": Float64Col(shape=(), dflt=0.0, pos=2)},
  "info1": {
    "name": StringCol(itemsize=10, shape=(), dflt='', pos=0),
    "value": Float64Col(shape=(), dflt=0.0, pos=1)},
  "color": EnumCol(enum=Enum({'blue': 2, 'green': 1, 'red': 0}), dflt='red',
    base=UInt32Atom(shape=(), dflt=0), shape=(), pos=2)}
```

As you can see, it provides very useful information on both the formats and the structure of the columns in your table.

This object also provides a natural naming approach to access to subcolumns metadata:

```
>>> table.description.info1
{"name": StringCol(itemsize=10, shape=(), dflt='', pos=0),
 "value": Float64Col(shape=(), dflt=0.0, pos=1)}
>>> table.description.info2.info3
{"x": Float64Col(shape=(), dflt=1.0, pos=0),
 "y": UInt8Col(shape=(), dflt=1, pos=1)}
```

There are other variables that can be interesting for you:

```
>>> table.description._v_nested_names
[('info2', [(('info3', [('x', 'y'])), 'name', 'value']),
 ('info1', ['name', 'value']), 'color']
>>> table.description.info1._v_nested_names
['name', 'value']
```

`_v_nested_names` provides the names of the columns as well as its structure. You can see that there are the same attributes for the different levels of the Description object, because the levels are *also* Description objects themselves.

There is a special attribute, called `_v_nested_descr`, that can be useful to create nested structured arrays that imitate the structure of the table (or a subtable thereof):

```
>>> import numpy
>>> table.description._v_nested_descr
[('info2', [(('info3', [('x', '()f8'), ('y', '()u1')]), ('name', '()S10'),
 ('value', '()f8')]), ('info1', [(('name', '()S10'), ('value', '()f8')]),
 ('color', '()u4')])
>>> numpy.rec.array(None, shape=0,
                    dtype=table.description._v_nested_descr)
recarray([],
          dtype=[('info2', [(('info3', [('x', '>f8'), ('y', '|u1')]),
 ('name', '|S10'), ('value', '>f8')]),
 ('info1', [(('name', '|S10'), ('value', '>f8')]),
 ('color', '>u4')])
>>> numpy.rec.array(None, shape=0,
                    dtype=table.description.info2._v_nested_descr)
recarray([],
          dtype=[('info3', [('x', '>f8'), ('y', '|u1')]), ('name', '|S10'),
 ('value', '>f8')])
```

You can see a simple example on how to create an array with NumPy.

Finally, there is a special iterator of the Description class, called `_f_walk` that is able to return you the different columns of the table:

```
>>> for coldescr in table.description._f_walk():
...     print("column-->", coldescr)
column--> Description([('info2', [(('info3', [('x', '()f8'), ('y', '()u1')]),
 ('name', '()S10'), ('value', '()f8')]),
 ('info1', [(('name', '()S10'), ('value', '()f8')]),
 ('color', '()u4')])
column--> EnumCol(enum=Enum({'blue': 2, 'green': 1, 'red': 0}), dflt='red',
                 base=UInt32Atom(shape=(), dflt=0), shape=(), pos=2)
column--> Description([('info3', [('x', '()f8'), ('y', '()u1')]), ('name', '()S10'),
 ('value', '()f8')])
column--> StringCol(itemsize=10, shape=(), dflt='', pos=1)
column--> Float64Col(shape=(), dflt=0.0, pos=2)
column--> Description([('name', '()S10'), ('value', '()f8')])
column--> StringCol(itemsize=10, shape=(), dflt='', pos=0)
column--> Float64Col(shape=(), dflt=0.0, pos=1)
column--> Description([('x', '()f8'), ('y', '()u1')])
column--> Float64Col(shape=(), dflt=1.0, pos=0)
column--> UInt8Col(shape=(), dflt=1, pos=1)
```

See the [The Description class](#) for the complete listing of attributes and methods of Description.

Well, this is the end of this tutorial. As always, do not forget to close your files:

```
>>> fileh.close()
```

Finally, you may want to have a look at your resulting data file.

```
$ ptdump -d nested-tut.h5
/ (RootGroup) ''
/table (Table(13,)) ''
  Data dump:
[0] (((1.0, 0), 'name2-0', 0.0), ('name1-0', 0.0), 0L)
[1] (((1.0, 1), 'name2-1', 0.0), ('name1-1', 0.0), 1L)
[2] (((1.0, 2), 'name2-2', 0.0), ('name1-2', 0.0), 2L)
[3] (((1.0, 3), 'name2-3', 0.0), ('name1-3', 0.0), 0L)
[4] (((1.0, 4), 'name2-4', 0.0), ('name1-4', 0.0), 1L)
[5] (((1.0, 5), 'name2-5', 0.0), ('name1-5', 0.0), 2L)
[6] (((1.0, 6), 'name2-6', 0.0), ('name1-6', 0.0), 0L)
[7] (((1.0, 7), 'name2-7', 0.0), ('name1-7', 0.0), 1L)
[8] (((1.0, 8), 'name2-8', 0.0), ('name1-8', 0.0), 2L)
[9] (((1.0, 9), 'name2-9', 0.0), ('name1-9', 0.0), 0L)
[10] (((1.0, 0), 'name2-0', 0.0), ('name1-0', 0.0), 0L)
[11] (((1.0, 4), 'name2-4', 0.0), ('name1-4', 0.0), 1L)
[12] (((1.0, 8), 'name2-8', 0.0), ('name1-8', 0.0), 2L)
/table2 (Table(3,)) ''
  Data dump:
[0] (((1.0, 0), 'name2-0', 0.0), ('name1-0', 0.0), 0L)
[1] (((1.0, 4), 'name2-4', 0.0), ('name1-4', 0.0), 1L)
[2] (((1.0, 8), 'name2-8', 0.0), ('name1-8', 0.0), 2L)
```

Most of the code in this section is also available in `examples/nested-tut.py`.

All in all, PyTables provides a quite comprehensive toolset to cope with nested structures and address your classification needs. However, caveat emptor, be sure to not nest your data too deeply or you will get inevitably messed interpreting too intertwined lists, tuples and description objects.

1.3.9 Other examples in PyTables distribution

Feel free to examine the rest of examples in directory `examples/`, and try to understand them. We have written several practical sample scripts to give you an idea of the PyTables capabilities, its way of dealing with HDF5 objects, and how it can be used in the real world.

1.4 Library Reference

PyTables implements several classes to represent the different nodes in the object tree. They are named `File`, `Group`, `Leaf`, `Table`, `Array`, `CArray`, `EArray`, `VArray` and `UnImplemented`. Another one allows the user to complement the information on these different objects; its name is `AttributeSet`. Finally, another important class called `IsDescription` allows to build a `Table` record description by declaring a subclass of it. Many other classes are defined in PyTables, but they can be regarded as helpers whose goal is mainly to declare the *data type properties* of the different first class objects and will be described at the end of this chapter as well.

An important function, called `open_file` 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 `is_pytables_file()` and `is_hdf5_file()`, respectively. There exists also a function called `which_lib_version()` that informs about the versions of the underlying C libraries (for example, *HDF5* or *Zlib*) and another called `print_versions()` that prints all the versions

of the software that PyTables relies on. Finally, `test()` lets you run the complete test suite from a Python console interactively.

1.4.1 Top-level variables and functions

Global variables

Global functions

1.4.2 File manipulation class

The File Class

File properties

File methods - file handling

File methods - hierarchy manipulation

File methods - tree traversal

File methods - Undo/Redo support

File methods - attribute handling

1.4.3 Hierarchy definition classes

The Node class

Node instance variables - location dependent

Node instance variables - location independent

Node instance variables - attribute shorthands

Node methods - hierarchy manipulation

Node methods - attribute handling

The Group class

Group properties

Group methods

Important: *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 (see *The File Class*, because they are most used in documentation and examples, and are a bit more powerful than those exposed here.

The following methods are provided in addition to those in Node (see *The Node class*):

Group special methods

Following are described the methods that automatically trigger actions when a Group instance is accessed in a special way.

This class defines the `__setattr__()`, `__getattr__()` and `__delattr__()` methods, and they set, get and delete *ordinary Python attributes* as normally intended. In addition to that, `__getattr__()` allows getting *child nodes* by their name for the sake of easy interaction on the command line, as long as there is no Python attribute with the same name. Groups also allow the interactive completion (when using readline) of the names of child nodes. For instance:

```
# get a Python attribute
nchild = group._v_nchildren

# Add a Table child called 'table' under 'group'.
h5file.create_table(group, 'table', my_description)
table = group.table           # get the table child instance
group.table = 'foo'          # set a Python attribute

# (PyTables warns you here about using the name of a child node.)
foo = group.table             # get a Python attribute
del group.table               # delete a Python attribute
table = group.table           # get the table child instance again
```

The Leaf class

Leaf properties

Leaf.size_in_memory

The size of this leaf's data in bytes when it is fully loaded into memory.

Leaf instance variables - aliases

The following are just easier-to-write aliases to their Node (see *The Node class*) counterparts (indicated between parentheses):

Leaf methods

1.4.4 Structured storage classes

The Table class

Table properties

Table methods - reading

Table methods - writing

Table methods - querying

Table methods - other

The Description class

Description methods

The Row class

Row methods

Row special methods

The Cols class

Cols properties

Cols methods

The Column class

Column instance variables

Column methods

Column special methods

1.4.5 Homogenous storage classes

The Array class

Array instance variables

`Array.atom`

An Atom (see *The Atom class and its descendants*) instance representing the *type* and *shape* of the atomic objects to be saved.

`Array.nrow`

On iterators, this is the index of the current row.

Array methods

Array special methods

The following methods automatically trigger actions when an `Array` instance is accessed in a special way (e.g. `array[2:3,...,:2]` will be equivalent to a call to `array.__getitem__((slice(2, 3, None), Ellipsis, slice(None, None, 2)))`).

The `CArray` class

The `EArray` class

`EArray` methods

The `VArray` class

`VArray` properties

`VArray` methods

`VArray` special methods

The following methods automatically trigger actions when a `VArray` instance is accessed in a special way (e.g., `varray[2:5]` will be equivalent to a call to `varray.__getitem__(slice(2, 5, None))`).

1.4.6 Link classes

The `Link` class

Link instance variables

Link methods

The following methods are useful for copying, moving, renaming and removing links.

The `SoftLink` class

`SoftLink` special methods

The following methods are specific for dereferencing and representing soft links.

The `ExternalLink` class

`ExternalLink` methods

ExternalLink special methods

The following methods are specific for dereferencing and representing external links.

1.4.7 Declarative classes

In this section a series of classes that are meant to *declare* datatypes that are required for creating primary PyTables datasets are described.

The Atom class and its descendants

Atom properties

Atom methods

Atom factory methods

Atom Sub-classes

Pseudo atoms

Now, there come three special classes, ObjectAtom, VLStringAtom and VLUnicodeAtom, that actually do not descend from Atom, but which goal is so similar that they should be described here. Pseudo-atoms can only be used with VLArray datasets (see *The VLArray class*), and they do not support multidimensional values, nor multiple values per row.

They can be recognised because they also have kind, type and shape attributes, but no size, itemsize or dtype ones. Instead, they have a base atom which defines the elements used for storage.

See `examples/vlarray1.py` and `examples/vlarray2.py` for further examples on VLArray datasets, including object serialization and string management.

ObjectAtom

VLStringAtom

VLUnicodeAtom

The Col class and its descendants

Col instance variables

In addition to the variables that they inherit from the Atom class, Col instances have the following attributes.

`Col._v_pos`

The *relative* position of this column with regard to its column siblings.

Col factory methods

Col sub-classes

The IsDescription class

Description helper functions

The AttributeSet class

AttributeSet properties

AttributeSet methods

1.4.8 Helper classes

This section describes some classes that do not fit in any other section and that mainly serve for ancillary purposes.

The Filters class

Filters methods

The Index class

Index instance variables

`tables.index.Index.nelements`

The number of currently indexed rows for this column.

Index methods

Index special methods

The IndexArray class

The Enum class

Enum special methods

The Unimplemented class

The Unknown class

Exceptions module

In the `exceptions` module exceptions and warnings that are specific to PyTables are declared.

1.4.9 General purpose expression evaluator class

The Expr class

Expr methods

Expr special methods

1.4.10 Filenode Module

Module constants

Module functions

The RawPyTablesIO base class

RawPyTablesIO attributes

RawPyTablesIO methods

The ROFileNode class

ROFileNode attributes

ROFileNode methods

The RAFileNode class

RAFileNode attributes

RAFileNode methods

1.5 Optimization tips

... durch planmässiges Tattonieren.

[... through systematic, palpable experimentation.]

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

On this chapter, you will get deeper knowledge of PyTables internals. PyTables has many tunable features so that you can improve the performance of your 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 efficiency boost for your code. But if your datasets are small (say, up to 10 MB) or your number of nodes is contained (up to 1000), you should not worry about that as the default parameters in PyTables are already tuned for those sizes (although you may want to adjust them further anyway). At any rate, reading this chapter will help you in your life with PyTables.

1.5.1 Understanding chunking

The underlying HDF5 library that is used by PyTables allows for certain datasets (the so-called *chunked* datasets) to take the data in bunches of a certain length, named *chunks*, and 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, check-summing, etc. on entire chunks.

HDF5 keeps a B-tree in memory that 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 cause 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 data (big B-trees).

In the next couple of sections, you will discover how to inform PyTables about the expected size of your datasets for allowing a sensible computation of the chunk sizes. Also, you will be presented some experiments so that you can get a feeling on the consequences of manually specifying the chunk size. Although doing this latter is only reserved to experienced people, these benchmarks may allow you to understand more deeply the chunk size implications and let you quickly start with the fine-tuning of this important parameter.

Informing PyTables about expected number of rows in tables or arrays

PyTables can determine a sensible chunk size to your dataset size if you help it by providing an estimation of the final number of rows for an extensible leaf¹. You should provide this information at leaf creation time by passing this value to the `expectedrows` argument of the `File.create_table()` method or `File.create_earray()` method (see [The EArray class](#)).

When your leaf size is bigger than 10 MB (take this figure only as a reference, not strictly), by providing this guess you will be optimizing the access to your data. When the table or array size is larger than, say 100MB, you are *strongly* suggested to provide such a guess; failing to do that may cause your application to do very slow I/O operations and to demand *huge* amounts of memory. You have been warned!

Fine-tuning the chunksize

Warning: This section is mostly meant for experts. If you are a beginner, you must know that setting manually the chunksize is a potentially dangerous action.

Most of the time, informing PyTables about the extent of your dataset is enough. However, for more sophisticated applications, when one has special requirements for doing the I/O or when dealing with really large datasets, you should really understand the implications of the chunk size in order to be able to find the best value for your own application.

You can specify the chunksize for every chunked dataset in PyTables by passing the `chunkshape` argument to the corresponding constructors. It is important to point out that `chunkshape` is not exactly the same thing than a `chunksize`; in fact, the `chunksize` of a dataset can be computed multiplying all the dimensions of the `chunkshape` among them and multiplying the outcome by the size of the atom.

We are going to describe a series of experiments where an EArray of 15 GB is written with different chunksizes, and then it is accessed in both sequential (i.e. first element 0, then element 1 and so on and so forth until the data is exhausted) and random mode (i.e. single elements are read randomly all through the dataset). These benchmarks have been carried out with PyTables 2.1 on a machine with an Intel Core2 processor @ 3 GHz and a RAID-0 made of two SATA disks spinning at 7200 RPM, and using GNU/Linux with an XFS filesystem. The script used for the benchmarks is available in `bench/optimal-chunksize.py`.

¹ CArray nodes, though not extensible, are chunked and have their optimum chunk size automatically computed at creation time, since their final shape is known.

In figures [Figure 1](#), [Figure 2](#), [Figure 3](#) and [Figure 4](#), you can see how the chunksize affects different aspects, like creation time, file sizes, sequential read time and random read time. So, if you properly inform PyTables about the extent of your datasets, you will get an automatic chunksize value (256 KB in this case) that is pretty optimal for most of uses. However, if what you want is, for example, optimize the creation time when using the Zlib compressor, you may want to reduce the chunksize to 32 KB (see [Figure 1](#)). Or, if your goal is to optimize the sequential access time for an dataset compressed with Blosc, you may want to increase the chunksize to 512 KB (see [Figure 3](#)).

You will notice that, by manually specifying the chunksize of a leave you will not normally get a drastic increase in performance, but at least, you have the opportunity to fine-tune such an important parameter for improve performance.

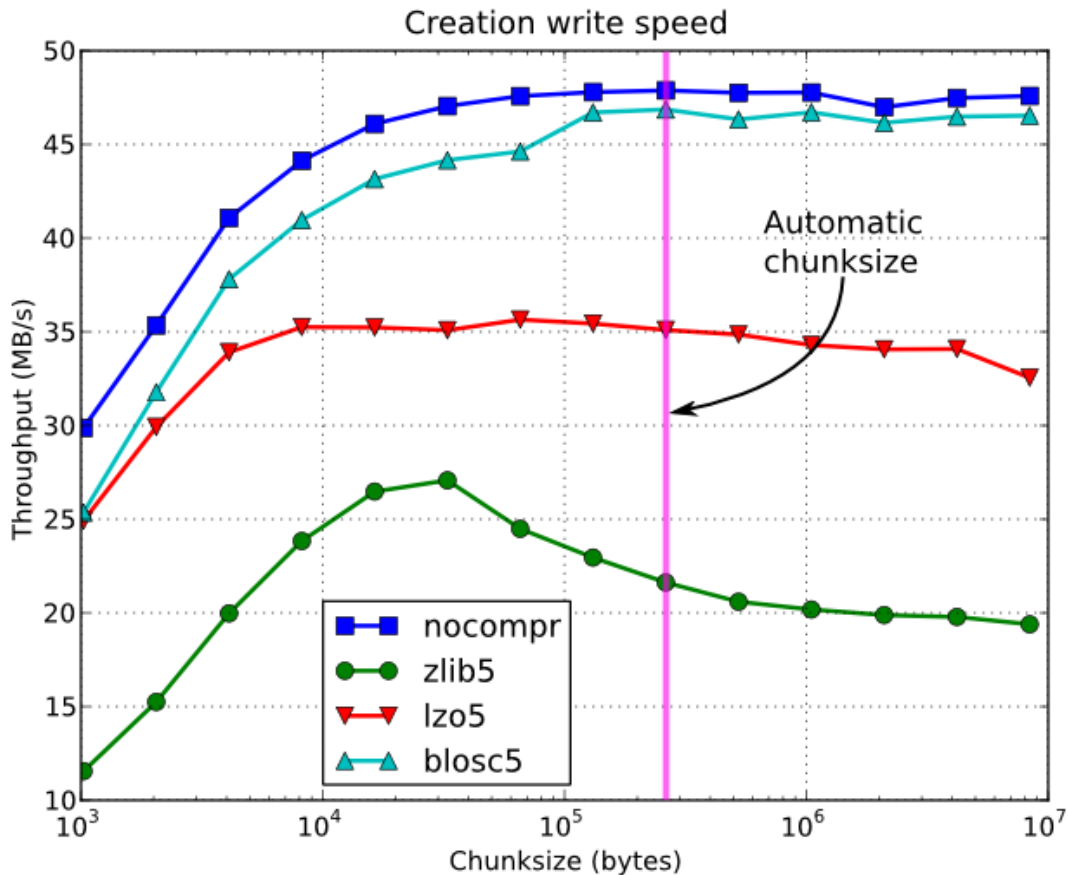


Fig. 7: **Figure 1.** Creation time per element for a 15 GB EArray and different chunksizes.

Finally, it is worth noting that adjusting the chunksize can be specially important if you want to access your dataset by blocks of certain dimensions. In this case, it is normally a good idea to set your chunkshape to be the same than these dimensions; you only have to be careful to not end with a too small or too large chunksize. As always, experimenting prior to pass your application into production is your best ally.

1.5.2 Accelerating your searches

Note: Many of the explanations and plots in this section and the forthcoming ones still need to be updated to include Blosc (see [BLOSC](#)), the new and powerful compressor added in PyTables 2.2 series. You should expect it to be the fastest compressor among all the described here, and its use is strongly recommended whenever you need extreme

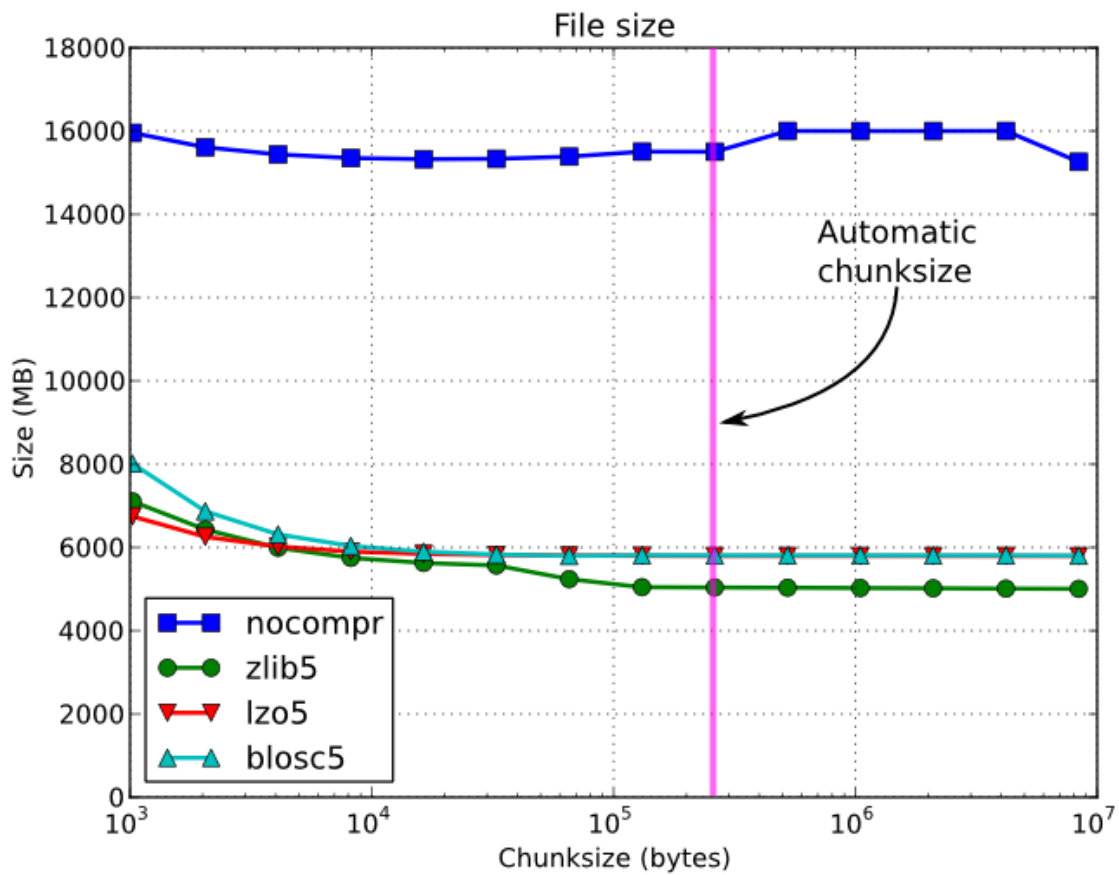


Fig. 8: Figure 2. File sizes for a 15 GB EArray and different chunksizes.

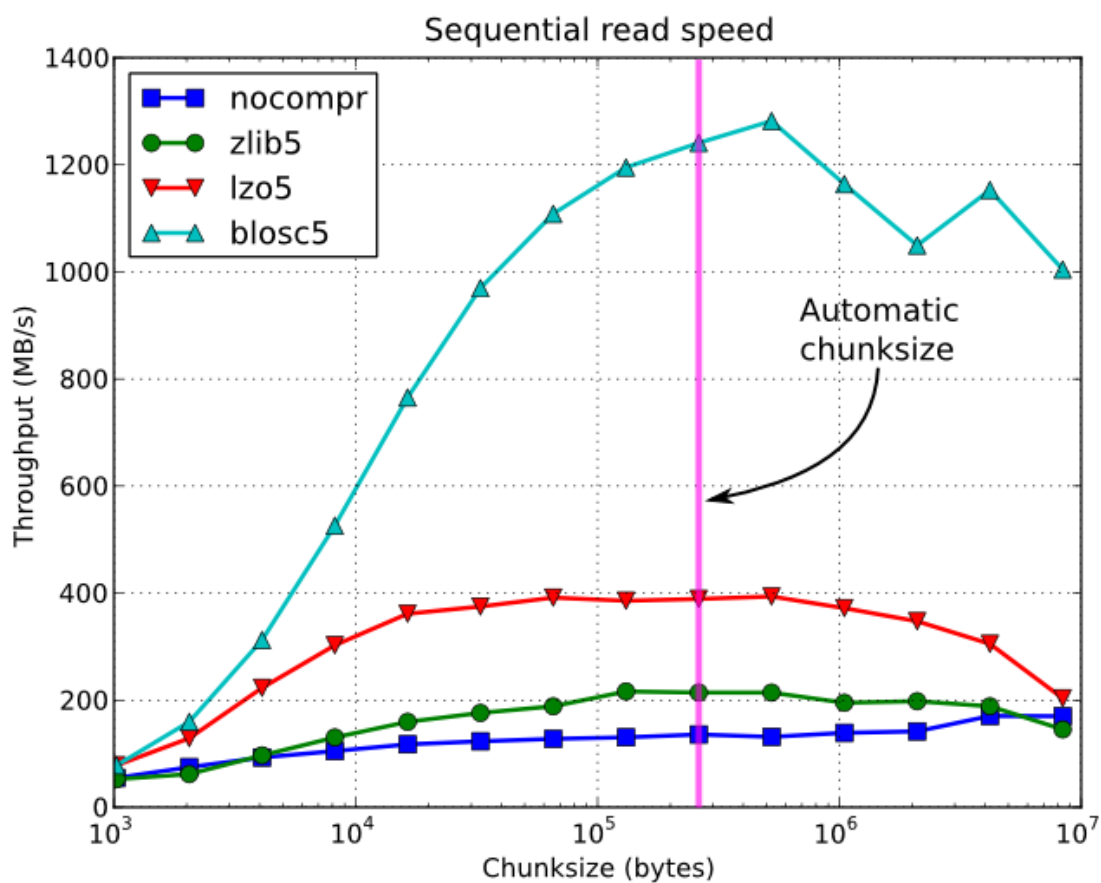


Fig. 9: Figure 3. Sequential access time per element for a 15 GB EArray and different chunksizes.

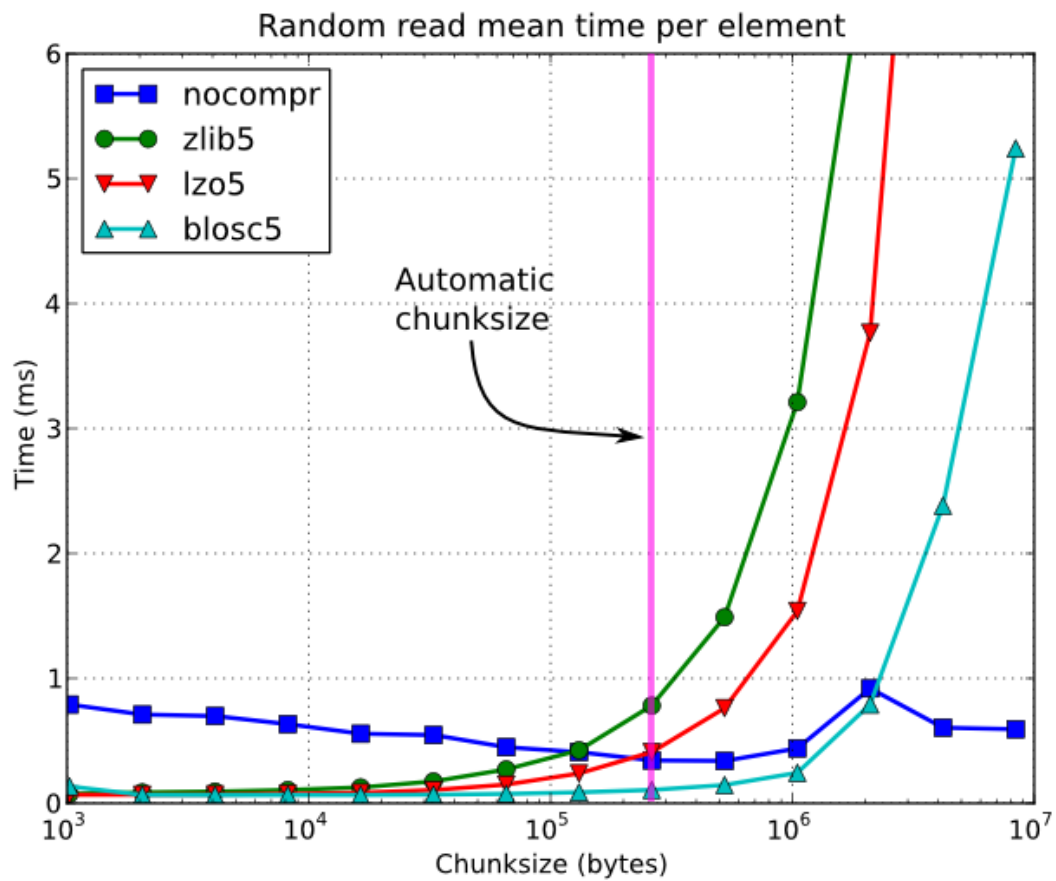


Fig. 10: **Figure 4.** Random access time per element for a 15 GB EArray and different chunksizes.

speed and not a very high compression ratio.

Searching in tables is one of the most common and time consuming operations that a typical user faces in the process of mining through his data. Being able to perform queries as fast as possible will allow more opportunities for finding the desired information quicker and also allows to deal with larger datasets.

PyTables offers many sort of techniques so as to speed-up the search process as much as possible and, in order to give you hints to use them based, a series of benchmarks have been designed and carried out. All the results presented in this section have been obtained with synthetic, random data and using PyTables 2.1. Also, the tests have been conducted on a machine with an Intel Core2 (64-bit) @ 3 GHz processor with RAID-0 disk storage (made of four spinning disks @ 7200 RPM), using GNU/Linux with an XFS filesystem. The script used for the benchmarks is available in `bench/indexed_search.py`. As your data, queries and platform may be totally different for your case, take this just as a guide because your mileage may vary (and will vary).

In order to be able to play with tables with a number of rows as large as possible, the record size has been chosen to be rather small (24 bytes). Here it is its definition:

```
class Record(tables.IsDescription):
    col1 = tables.Int32Col()
    col2 = tables.Int32Col()
    col3 = tables.Float64Col()
    col4 = tables.Float64Col()
```

In the next sections, we will be optimizing the times for a relatively complex query like this:

```
result = [row['col2'] for row in table if (
    ((row['col4'] >= lim1 and row['col4'] < lim2) or
    ((row['col2'] > lim3 and row['col2'] < lim4))) and
    ((row['col1']+3.1*row['col2']+row['col3']*row['col4']) > lim5)
)]
```

(for future reference, we will call this sort of queries *regular* queries). So, if you want to see how to greatly improve the time taken to run queries like this, keep reading.

In-kernel searches

PyTables provides a way to accelerate data selections inside of a single table, through the use of the *Table methods - querying* iterator and related query methods. This mode of selecting data is called *in-kernel*. Let's see an example of an *in-kernel* query based on the *regular* one mentioned above:

```
result = [row['col2'] for row in table.where(
    '''((col4 >= lim1) & (col4 < lim2)) |
    ((col2 > lim3) & (col2 < lim4)) &
    ((col1+3.1*col2+col3*col4) > lim5)''')]
```

This simple change of mode selection can improve search times quite a lot and actually make PyTables very competitive when compared against typical relational databases as you can see in [Figure 5](#) and [Figure 6](#).

By looking at [Figure 5](#) you can see how in the case that table data fits easily in memory, in-kernel searches on uncompressed tables are generally much faster (10x) than standard queries as well as PostgreSQL (5x). Regarding compression, we can see how Zlib compressor actually slows down the performance of in-kernel queries by a factor 3.5x; however, it remains faster than PostgreSQL (40%). On his hand, LZO compressor only decreases the performance by a 75% with respect to uncompressed in-kernel queries and is still a lot faster than PostgreSQL (3x). Finally, one can observe that, for low selectivity queries (large number of hits), PostgreSQL performance degrades quite steadily, while in PyTables this slow down rate is significantly smaller. The reason of this behaviour is not entirely clear to the authors, but the fact is clearly reproducible in our benchmarks.

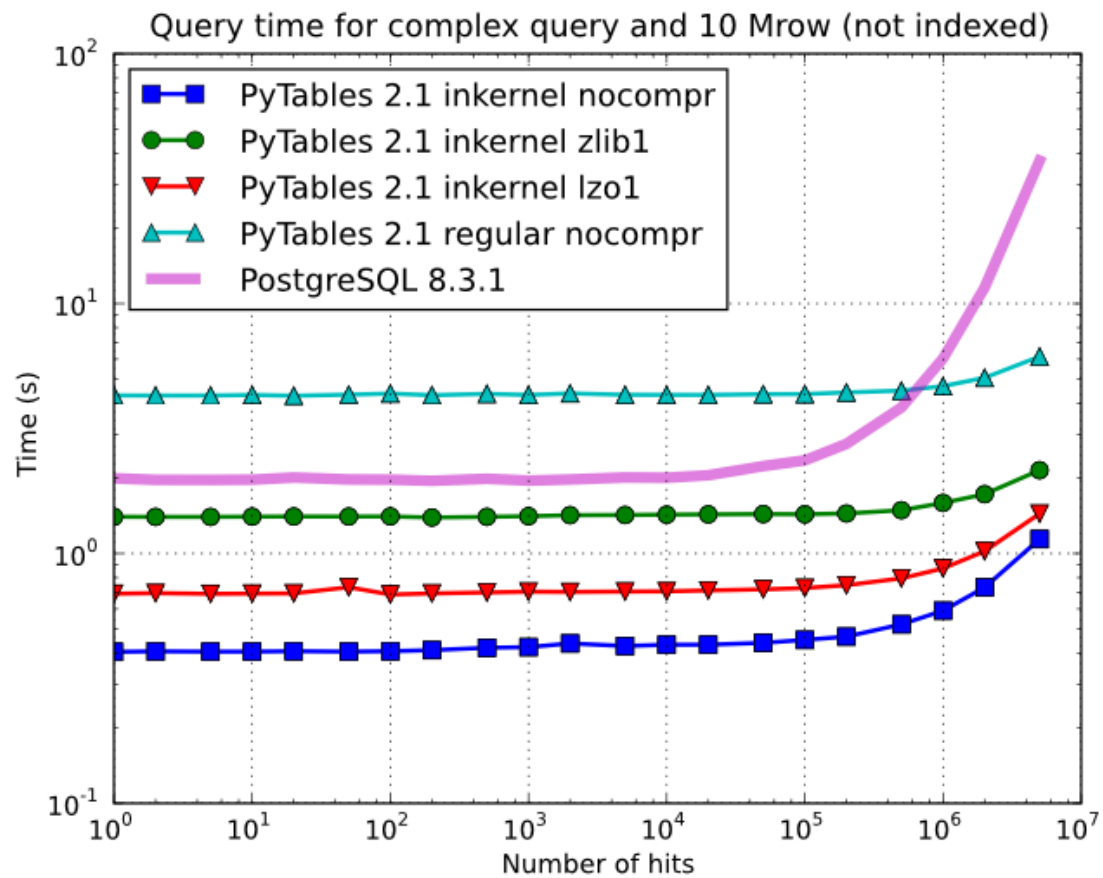


Fig. 11: Figure 5. Times for non-indexed complex queries in a small table with 10 millions of rows: the data fits in memory.

But, why in-kernel queries are so fast when compared with regular ones?. The answer is that in regular selection mode the data for all the rows in table has to be brought into Python space so as to evaluate the condition and decide if the corresponding field should be added to the result list. On the contrary, in the in-kernel mode, the condition is passed to the PyTables kernel (hence the name), written in C, and evaluated there at full C speed (with the help of the integrated Numexpr package, see [\[NUMEXPR\]](#)), so that the only values that are brought to Python space are the rows that fulfilled the condition. Hence, for selections that only have a relatively small number of hits (compared with the total amount of rows), the savings are very large. It is also interesting to note the fact that, although for queries with a large number of hits the speed-up is not as high, it is still very important.

On the other hand, when the table is too large to fit in memory (see [Figure 6](#)), the difference in speed between regular and in-kernel is not so important, but still significant (2x). Also, and curiously enough, large tables compressed with Zlib offers slightly better performance (around 20%) than uncompressed ones; this is because the additional CPU spent by the uncompressor is compensated by the savings in terms of net I/O (one has to read less actual data from disk). However, when using the extremely fast LZO compressor, it gives a clear advantage over Zlib, and is up to 2.5x faster than not using compression at all. The reason is that LZO decompression speed is much faster than Zlib, and that allows PyTables to read the data at full disk speed (i.e. the bottleneck is in the I/O subsystem, not in the CPU). In this case the compression rate is around 2.5x, and this is why the data can be read 2.5x faster. So, in general, using the LZO compressor is the best way to ensure best reading/querying performance for out-of-core datasets (more about how compression affects performance in [Compression issues](#)).

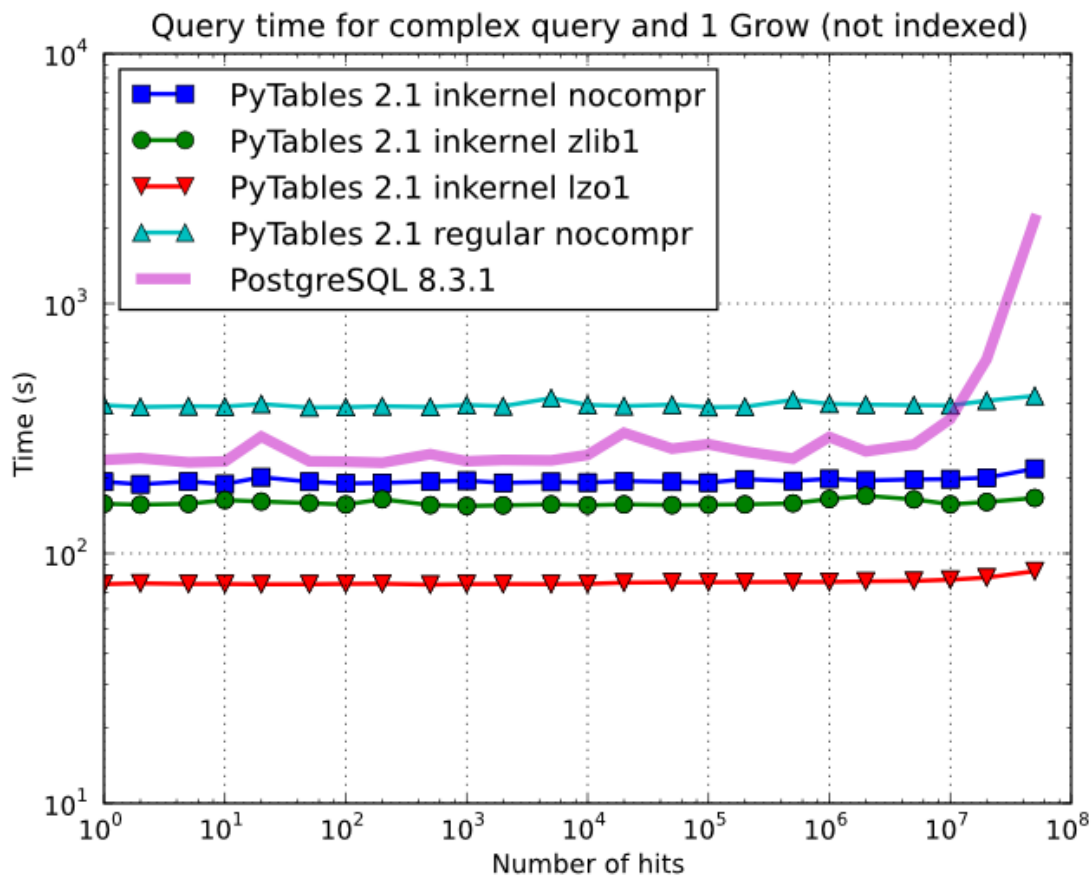


Fig. 12: Figure 6. Times for non-indexed complex queries in a large table with 1 billion of rows: the data does not fit in memory.

Furthermore, you can mix the *in-kernel* and *regular* selection modes for evaluating arbitrarily complex conditions

making use of external functions. Look at this example:

```
result = [ row['var2']
           for row in table.where('(var3 == "foo") & (var1 <= 20)')
           if your_function(row['var2']) ]
```

Here, we use an *in-kernel* selection to choose rows according to the values of the `var3` and `var1` fields. Then, we apply a *regular* selection to complete the query. Of course, when you mix the *in-kernel* and *regular* selection modes you should pass the most restrictive condition to the *in-kernel* part, i.e. to the `where()` iterator. In situations where it is not clear which is the most restrictive condition, you might want to experiment a bit in order to find the best combination.

However, since in-kernel condition strings allow rich expressions allowing the coexistence of multiple columns, variables, arithmetic operations and many typical functions, it is unlikely that you will be forced to use external regular selections in conditions of small to medium complexity. See [Condition Syntax](#) for more information on in-kernel condition syntax.

Indexed searches

When you need more speed than *in-kernel* selections can offer you, PyTables offers a third selection method, the so-called *indexed* mode (based on the highly efficient OPSI indexing engine). In this mode, you have to decide which column(s) you are going to apply your selections over, and index them. Indexing is just a kind of sorting operation over a column, so that searches along such a column (or columns) will look at this sorted information by using a *binary search* which is much faster than the *sequential search* described in the previous section.

You can index the columns you want by calling the `Column.create_index()` method on an already created table. For example:

```
indexrows = table.cols.var1.create_index()
indexrows = table.cols.var2.create_index()
indexrows = table.cols.var3.create_index()
```

will create indexes for all `var1`, `var2` and `var3` columns.

After you have indexed a series of columns, the PyTables query optimizer will try hard to discover the usable indexes in a potentially complex expression. However, there are still places where it cannot determine that an index can be used. See below for examples where the optimizer can safely determine if an index, or series of indexes, can be used or not.

Example conditions where an index can be used:

- `var1 >= "foo"` (`var1` is used)
- `var1 >= mystr` (`var1` is used)
- `(var1 >= "foo") & (var4 > 0.0)` (`var1` is used)
- `("bar" <= var1) & (var1 < "foo")` (`var1` is used)
- `((("bar" <= var1) & (var1 < "foo")) & (var4 > 0.0))` (`var1` is used)
- `(var1 >= "foo") & (var3 > 10)` (`var1` and `var3` are used)
- `(var1 >= "foo") | (var3 > 10)` (`var1` and `var3` are used)
- `~(var1 >= "foo") | ~(var3 > 10)` (`var1` and `var3` are used)

Example conditions where an index can *not* be used:

- `var4 > 0.0` (`var4` is not indexed)
- `var1 != 0.0` (range has two pieces)

- `~(("bar" <= var1) & (var1 < "foo")) & (var4 > 0.0)` (negation of a complex boolean expression)

Note: From PyTables 2.3 on, several indexes can be used in a single query.

Note: If you want to know for sure whether a particular query will use indexing or not (without actually running it), you are advised to use the `Table.will_query_use_indexing()` method.

One important aspect of the new indexing in PyTables (≥ 2.3) is that it has been designed from the ground up with the goal of being capable to effectively manage very large tables. To this goal, it sports a wide spectrum of different quality levels (also called optimization levels) for its indexes so that the user can choose the best one that suits her needs (more or less size, more or less performance).

In [Figure 7](#), you can see that the times to index columns in tables can be really short. In particular, the time to index a column with 1 billion rows (1 Gigarow) with the lowest optimization level is less than 4 minutes while indexing the same column with full optimization (so as to get a completely sorted index or CSI) requires around 1 hour. These are rather competitive figures compared with a relational database (in this case, PostgreSQL 8.3.1, which takes around 1.5 hours for getting the index done). This is because PyTables is geared towards read-only or append-only tables and takes advantage of this fact to optimize the indexes properly. On the contrary, most relational databases have to deliver decent performance in other scenarios as well (specially updates and deletions), and this fact leads not only to slower index creation times, but also to indexes taking much more space on disk, as you can see in [Figure 8](#).

The user can select the index quality by passing the desired `optlevel` and `kind` arguments to the `Column.create_index()` method. We can see in figures [Figure 7](#) and [Figure 8](#) how the different optimization levels affects index time creation and index sizes.

So, which is the effect of the different optimization levels in terms of query times? You can see that in [Figure 9](#).

Of course, compression also has an effect when doing indexed queries, although not very noticeable, as can be seen in [Figure 10](#). As you can see, the difference between using no compression and using Zlib or LZO is very little, although LZO achieves relatively better performance generally speaking.

You can find a more complete description and benchmarks about OPSI, the indexing system of PyTables (≥ 2.3) in [\[OPSI\]](#).

Indexing and Solid State Disks (SSD)

Lately, the long promised Solid State Disks (SSD for brevity) with decent capacities and affordable prices have finally hit the market and will probably stay in coexistence with the traditional spinning disks for the foreseeable future (separately or forming *hybrid* systems). SSD have many advantages over spinning disks, like much less power consumption and better throughput. But of paramount importance, specially in the context of accelerating indexed queries, is its very reduced latency during disk seeks, which is typically 100x better than traditional disks. Such a huge improvement has to have a clear impact in reducing the query times, specially when the selectivity is high (i.e. the number of hits is small).

In order to offer an estimate on the performance improvement we can expect when using a low-latency SSD instead of traditional spinning disks, the benchmark in the previous section has been repeated, but this time using a single SSD disk instead of the four spinning disks in RAID-0. The result can be seen in [Figure 11](#). There one can see how a query in a table of 1 billion of rows with 100 hits took just 1 tenth of second when using a SSD, instead of 1 second that needed the RAID made of spinning disks. This factor of 10x of speed-up for high-selectivity queries is nothing to sneeze at, and should be kept in mind when really high performance in queries is needed. It is also interesting that using compression with LZO does have a clear advantage over when no compression is done.

Finally, we should remark that SSD can't compete with traditional spinning disks in terms of capacity as they can only provide, for a similar cost, between 1/10th and 1/50th of the size of traditional disks. It is here where the compression

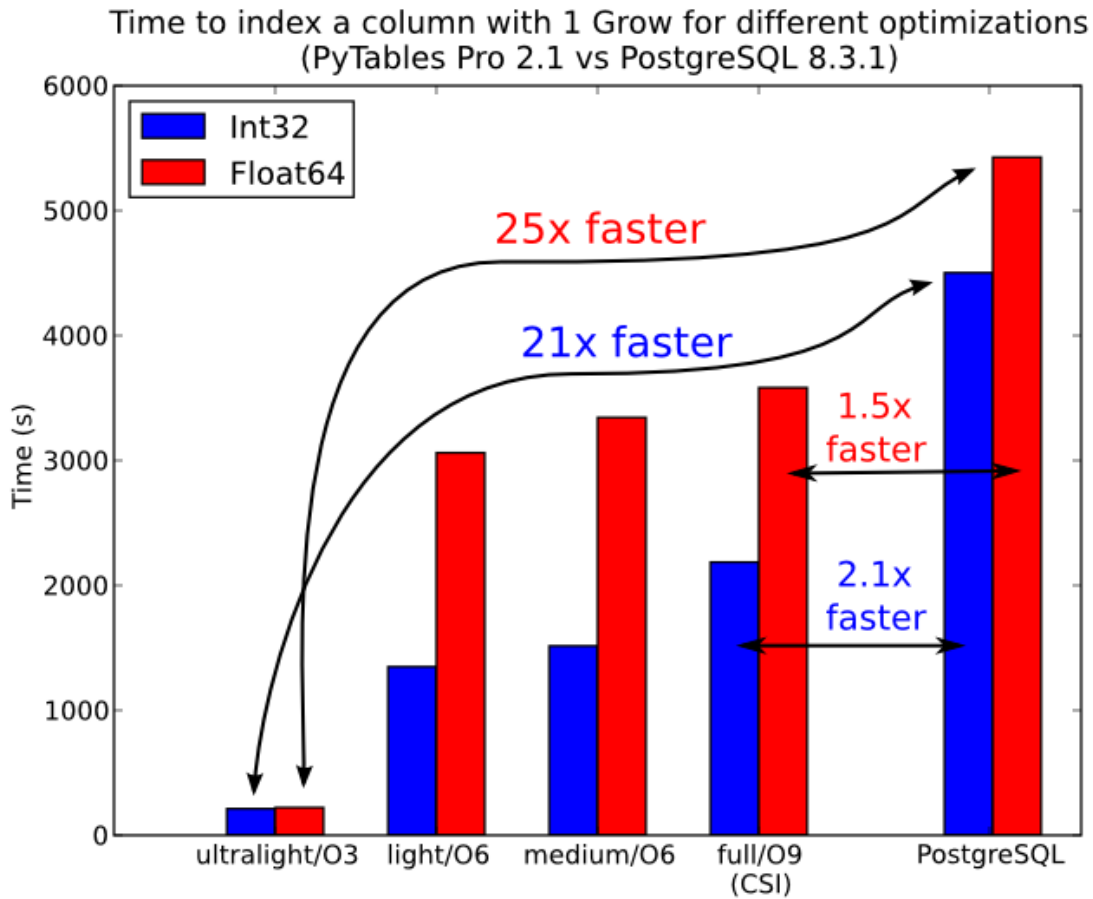


Fig. 13: Figure 7. Times for indexing an Int32 and Float64 column.

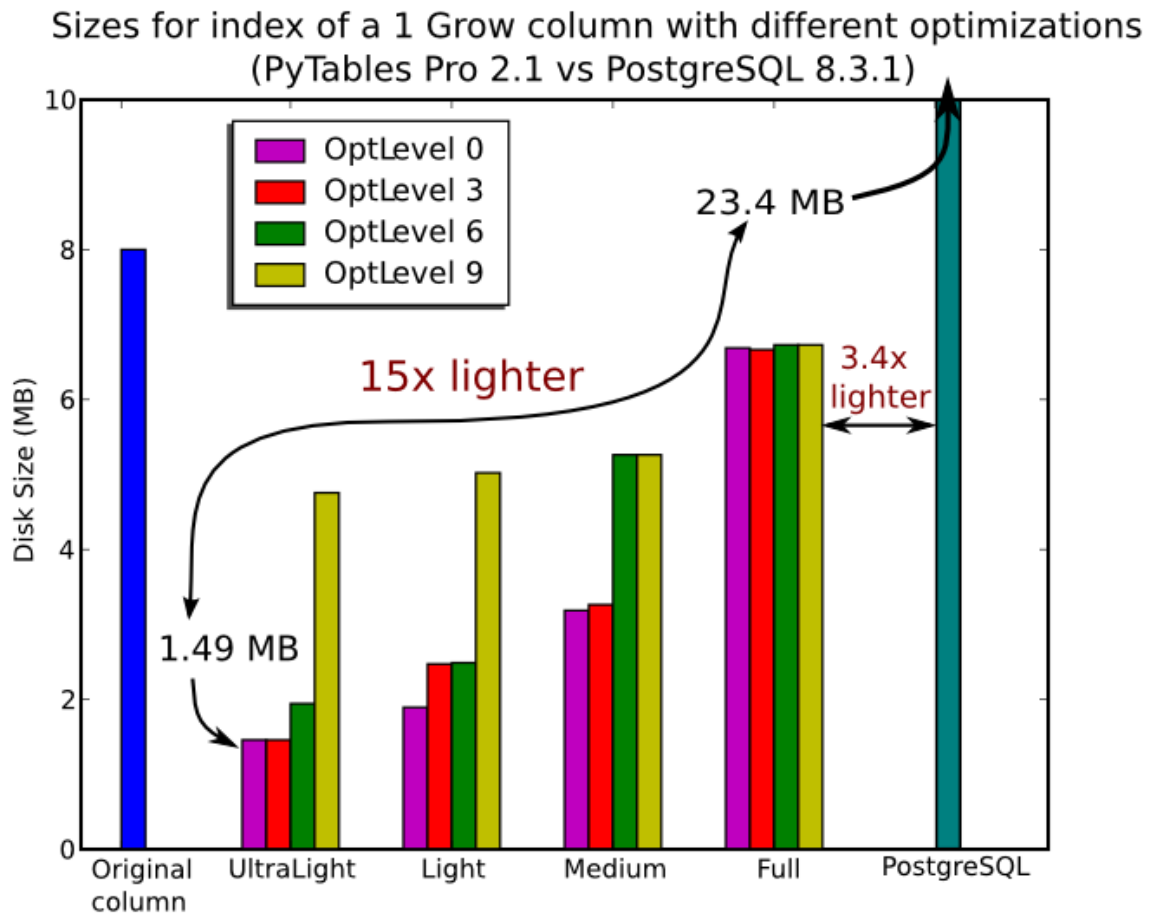


Fig. 14: Figure 8. Sizes for an index of a Float64 column with 1 billion of rows.

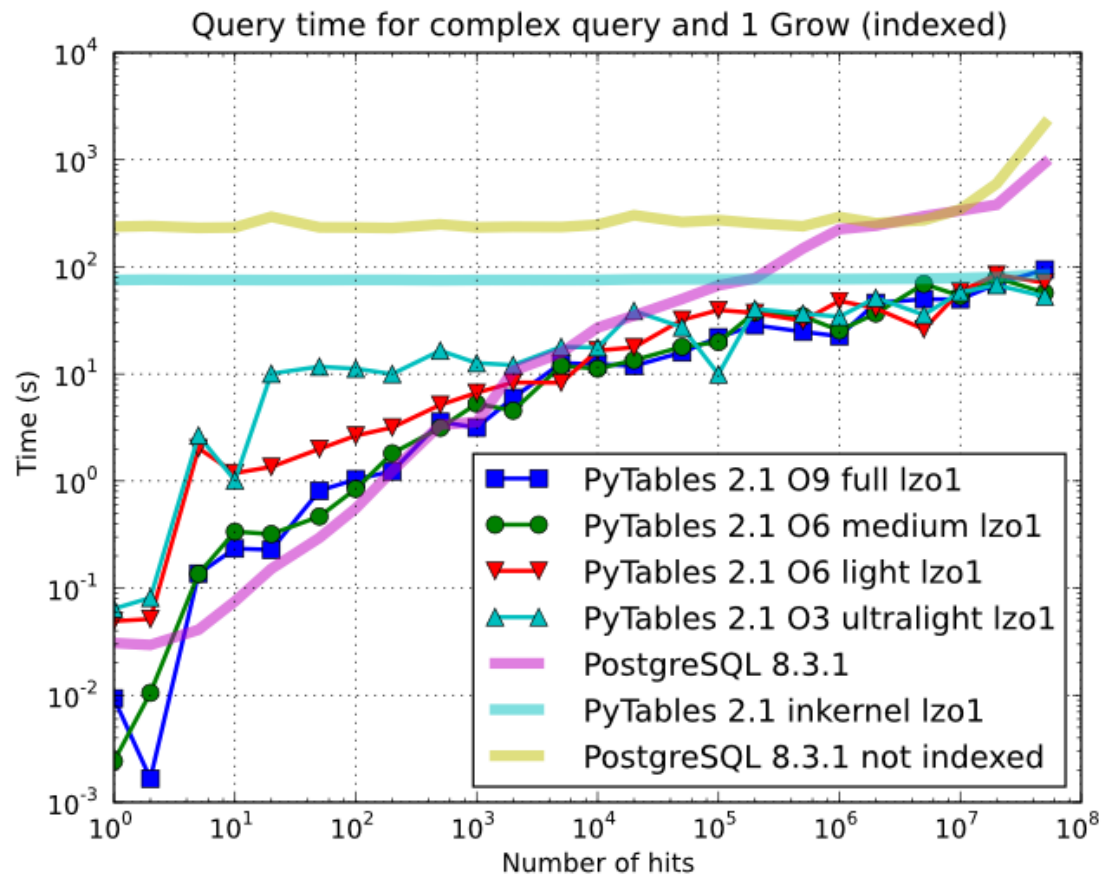


Fig. 15: Figure 9. Times for complex queries with a cold cache (mean of 5 first random queries) for different optimization levels. Benchmark made on a machine with Intel Core2 (64-bit) @ 3 GHz processor with RAID-0 disk storage.

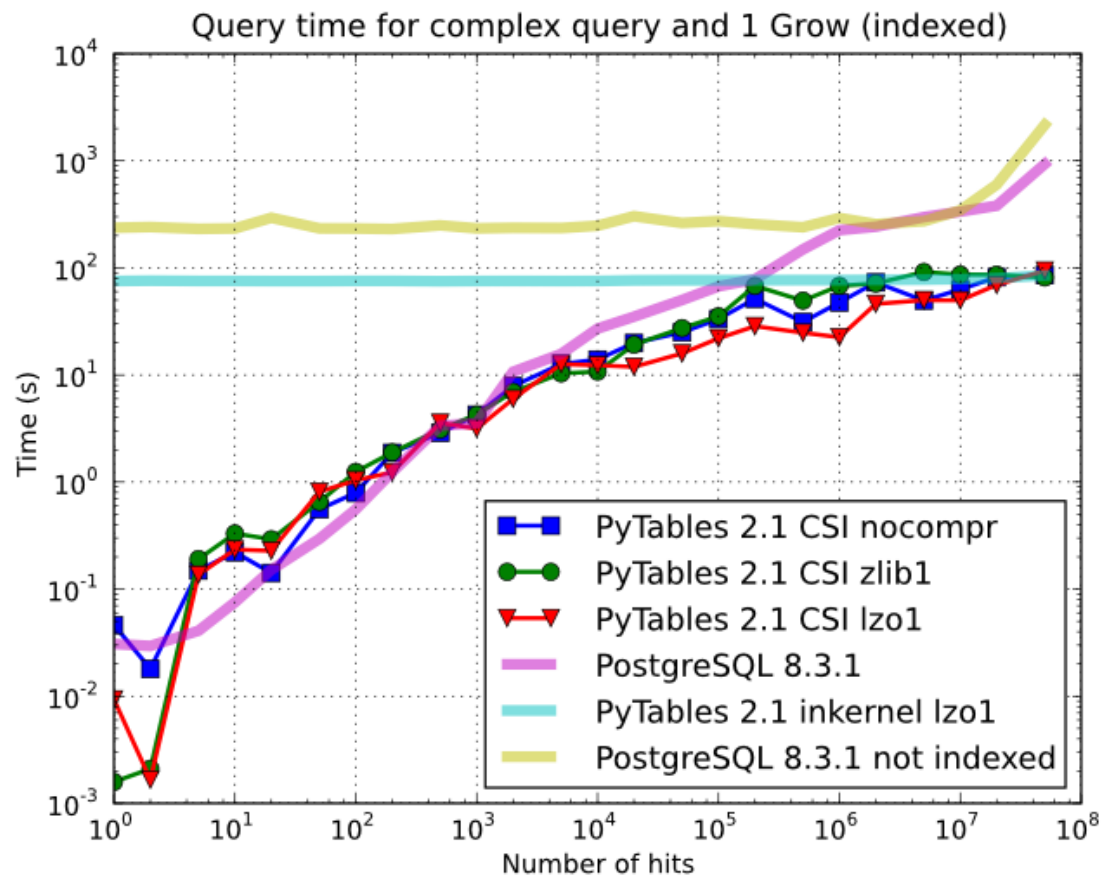


Fig. 16: Figure 10. Times for complex queries with a cold cache (mean of 5 first random queries) for different compressors.

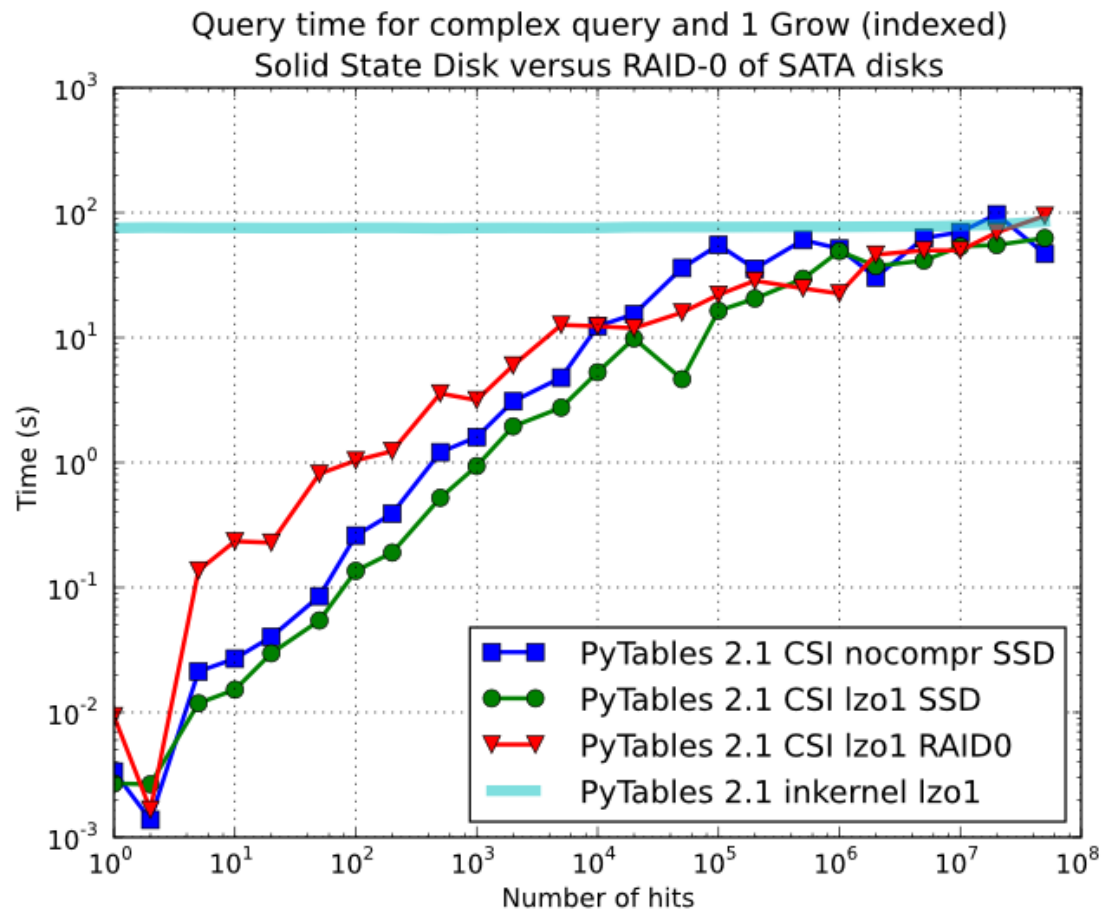


Fig. 17: Figure 11. Times for complex queries with a cold cache (mean of 5 first random queries) for different disk storage (SSD vs spinning disks).

capabilities of PyTables can be very helpful because both tables and indexes can be compressed and the final space can be reduced by typically 2x to 5x (4x to 10x when compared with traditional relational databases). Best of all, as already mentioned, performance is not degraded when compression is used, but actually *improved*. So, by using PyTables and SSD you can query larger datasets that otherwise would require spinning disks when using other databases

In fact, we were unable to run the PostgreSQL benchmark in this case because the space needed exceeded the capacity of our SSD., while allowing improvements in the speed of indexed queries between 2x (for medium to low selectivity queries) and 10x (for high selectivity queries).

Achieving ultimate speed: sorted tables and beyond

Warning: Sorting a large table is a costly operation. The next procedure should only be performed when your dataset is mainly read-only and meant to be queried many times.

When querying large tables, most of the query time is spent in locating the interesting rows to be read from disk. In some occasions, you may have queries whose result depends *mainly* of one single column (a query with only one single condition is the trivial example), so we can guess that sorting the table by this column would lead to locate the interesting rows in a much more efficient way (because they would be mostly *contiguous*). We are going to confirm this guess.

For the case of the query that we have been using in the previous sections:

```
result = [row['col2'] for row in table.where(
    '''((col4 >= lim1) & (col4 < lim2)) |
       ((col2 > lim3) & (col2 < lim4)) &
       ((col1+3.1*col2+col3*col4) > lim5))''')]
```

it is possible to determine, by analysing the data distribution and the query limits, that col4 is such a *main column*. So, by ordering the table by the col4 column (for example, by specifying setting the column to sort by in the `sortby` parameter in the `Table.copy()` method and re-indexing col2 and col4 afterwards, we should get much faster performance for our query. This is effectively demonstrated in [Figure 12](#), where one can see how queries with a low to medium (up to 10000) number of hits can be done in around 1 tenth of second for a RAID-0 setup and in around 1 hundredth of second for a SSD disk. This represents up to more that 100x improvement in speed with respect to the times with unsorted tables. On the other hand, when the number of hits is large (> 1 million), the query times grow almost linearly, showing a near-perfect scalability for both RAID-0 and SSD setups (the sequential access to disk becomes the bottleneck in this case).

Even though we have shown many ways to improve query times that should fulfill the needs of most of people, for those needing more, you can for sure discover new optimization opportunities. For example, querying against sorted tables is limited mainly by sequential access to data on disk and data compression capability, so you may want to read [Fine-tuning the chunksize](#), for ways on improving this aspect. Reading the other sections of this chapter will help in finding new roads for increasing the performance as well. You know, the limit for stopping the optimization process is basically your imagination (but, most plausibly, your available time ;-).

1.5.3 Compression issues

One of the beauties of PyTables is that it supports compression on tables and arrays², although it is not used by default. Compression of big amounts of data might be a bit controversial feature, because it has a legend of being a very big consumer of CPU time resources. However, if you are willing to check if compression can help not only by reducing your dataset file size but *also* by improving I/O efficiency, specially when dealing with very large datasets, keep reading.

² Except for Array objects.

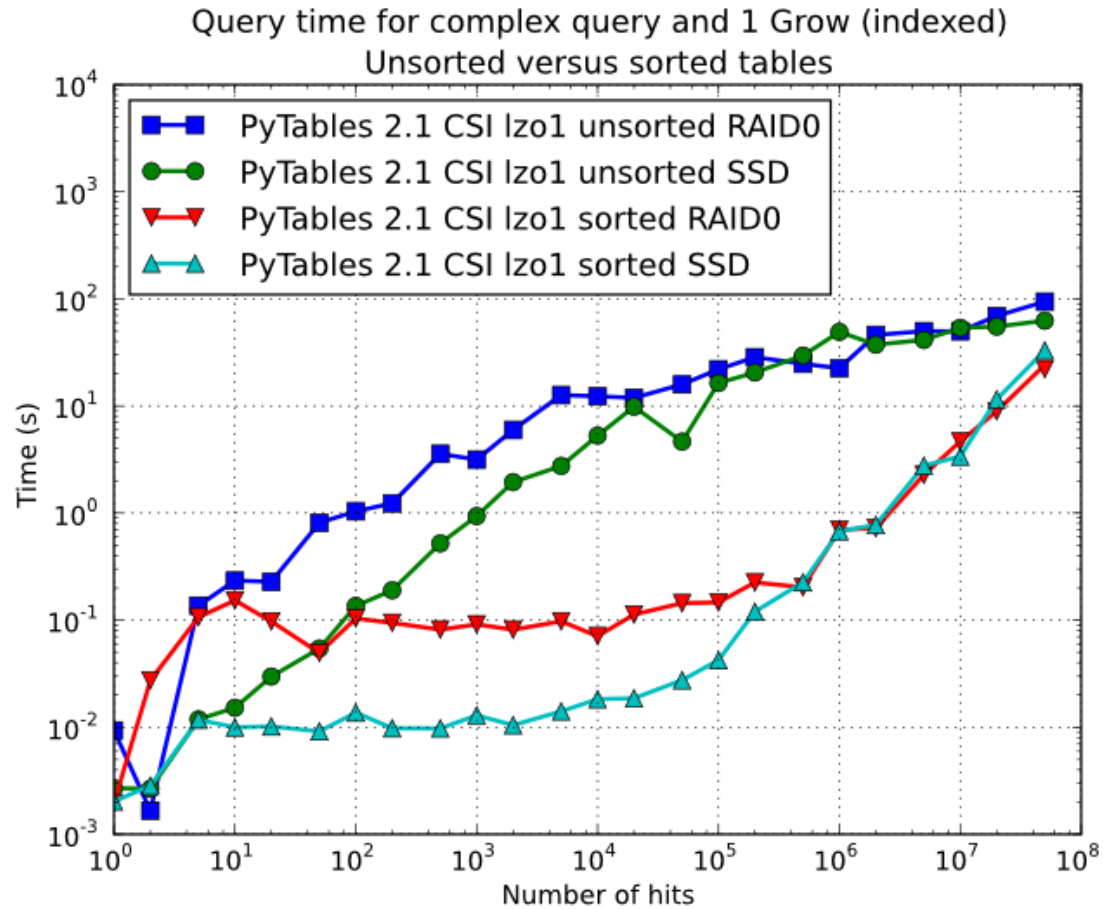


Fig. 18: **Figure 12.** Times for complex queries with a cold cache (mean of 5 first random queries) for unsorted and sorted tables.

A study on supported compression libraries

The compression library used by default is the *Zlib* (see [\[ZLIB\]](#)). Since *HDF5* *requires* it, you can safely use it and expect that your *HDF5* files will be readable 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 compressing you data.

However, in some situations it is critical to have a *very good decompression speed* (at the expense of lower compression ratios or more CPU wasted on compression, as we will see soon). In others, the emphasis is put in achieving the *maximum compression ratios*, no matter which reading speed will result. This is why support for two additional compressors has been added to *PyTables*: *LZO* (see [\[LZO\]](#)) and *bzip2* (see [\[BZIP2\]](#)). Following the author of *LZO* (and checked by the author of this section, as you will see soon), *LZO* offers pretty fast compression and extremely fast decompression. In fact, *LZO* is so fast when compressing/decompressing that it may well happen (that depends on your data, of course) that writing or reading a compressed dataset is sometimes faster than if it is not compressed at all (specially when dealing with extremely large datasets). This fact is very important, specially if you have to deal with very large amounts of data. Regarding *bzip2*, it has a reputation of achieving excellent compression ratios, but at the price of spending much more CPU time, which results in very low compression/decompression speeds.

Be aware that the *LZO* and *bzip2* 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. Still, see the [ptrepack](#) (where the *ptrepack* utility is described) to find a way to free your files from *LZO* or *bzip2* dependencies, so that you can use these compressors locally with the warranty that you can replace them with *Zlib* (or even remove compression completely) if you want to use these files with other *HDF5* tools or platforms afterwards.

In order to allow you to grasp what amount of compression can be achieved, and how this affects performance, a series of experiments has been carried out. All the results presented in this section (and in the next one) have been obtained with synthetic data and using *PyTables* 1.3. Also, the tests have been conducted on a IBM OpenPower 720 (e-series) with a PowerPC G5 at 1.65 GHz and a hard disk spinning at 15K RPM. As your data and platform may be totally different for your case, take this just as a guide because your mileage may vary. Finally, and to be able to play with tables with a number of rows as large as possible, the record size has been chosen to be small (16 bytes). Here is its definition:

```
class Bench(IsDescription):
    var1 = StringCol(length=4)
    var2 = IntCol()
    var3 = FloatCol()
```

With this setup, you can look at the compression ratios that can be achieved in [Figure 13](#). As you can see, *LZO* is the compressor that performs worse in this sense, but, curiously enough, there is not much difference between *Zlib* and *bzip2*.

Also, *PyTables* lets you select different compression levels for *Zlib* and *bzip2*, although you may get a bit disappointed by the small improvement that these compressors show when dealing with a combination of numbers and strings as in our example. As a reference, see plot [Figure 14](#) for a comparison of the compression achieved by selecting different levels of *Zlib*. Very oddly, the best compression ratio corresponds to level 1 (!). See later for an explanation and more figures on this subject.

Have also a look at [Figure 15](#). It shows how the speed of writing rows evolves as the size (number of rows) of the table grows. Even though in these graphs the size of one single row is 16 bytes, you can most probably extrapolate these figures to other row sizes.

In [Figure 16](#) you can see how compression affects the reading performance. In fact, what you see in the plot is an *in-kernel selection* speed, but provided that this operation is very fast (see [In-kernel searches](#)), we can accept it as an actual read test. Compared with the reference line without compression, the general trend here is that *LZO* does not affect too much the reading performance (and in some points it is actually better), *Zlib* makes speed drop to a half, while *bzip2* is performing very slow (up to 8x slower).

Also, in the same [Figure 16](#) you can notice some strange peaks in the speed that we might be tempted to attribute

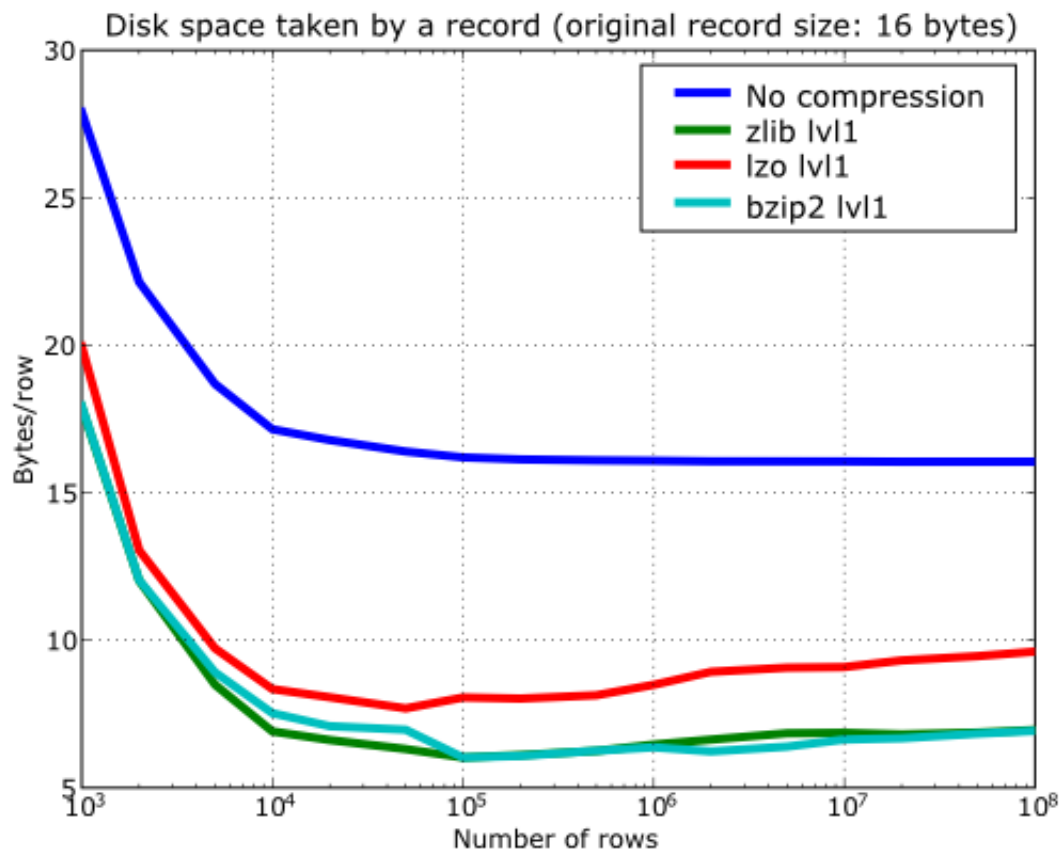


Fig. 19: Figure 13. Comparison between different compression libraries.

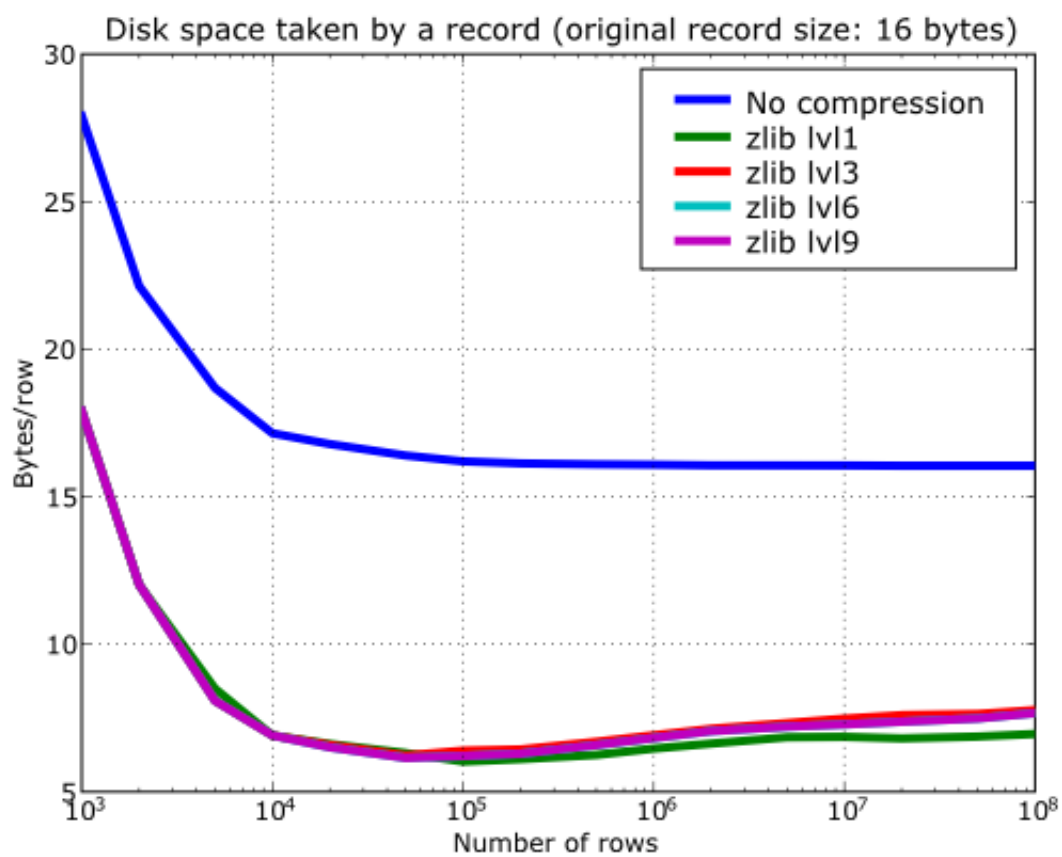


Fig. 20: **Figure 14.** Comparison between different compression levels of Zlib.

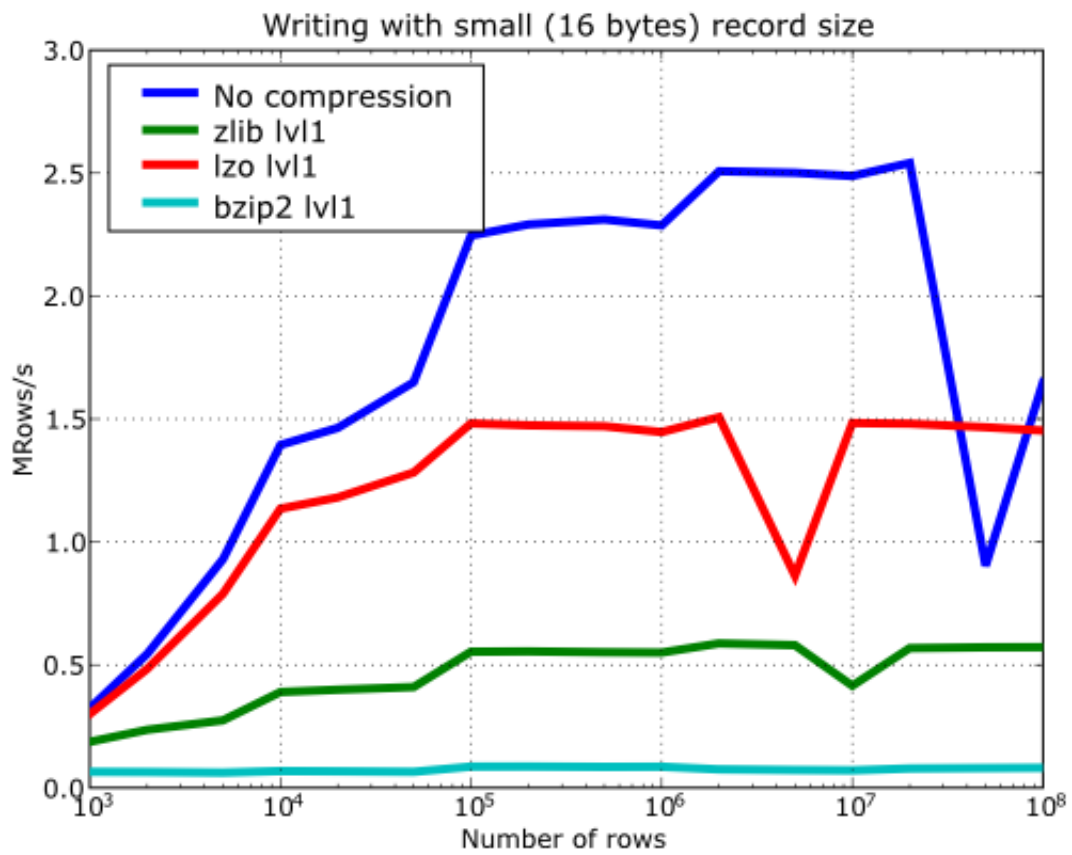


Fig. 21: Figure 15. Writing tables with several compressors.

to libraries on which PyTables relies (HDF5, compressors...), or to PyTables itself. However, [Figure 17](#) reveals that, if we put the file in the filesystem cache (by reading it several times before, for example), the evolution of the performance is much smoother. So, the most probable explanation would be that such peaks are a consequence of the underlying OS filesystem, rather than a flaw in PyTables (or any other library behind it). Another consequence that can be derived from the aforementioned plot is that LZO decompression performance is much better than Zlib, allowing an improvement in overall speed of more than 2x, and perhaps more important, the read performance for really large datasets (i.e. when they do not fit in the OS filesystem cache) can be actually *better* than not using compression at all. Finally, one can see that reading performance is very badly affected when bzip2 is used (it is 10x slower than LZO and 4x than Zlib), but this was somewhat expected anyway.

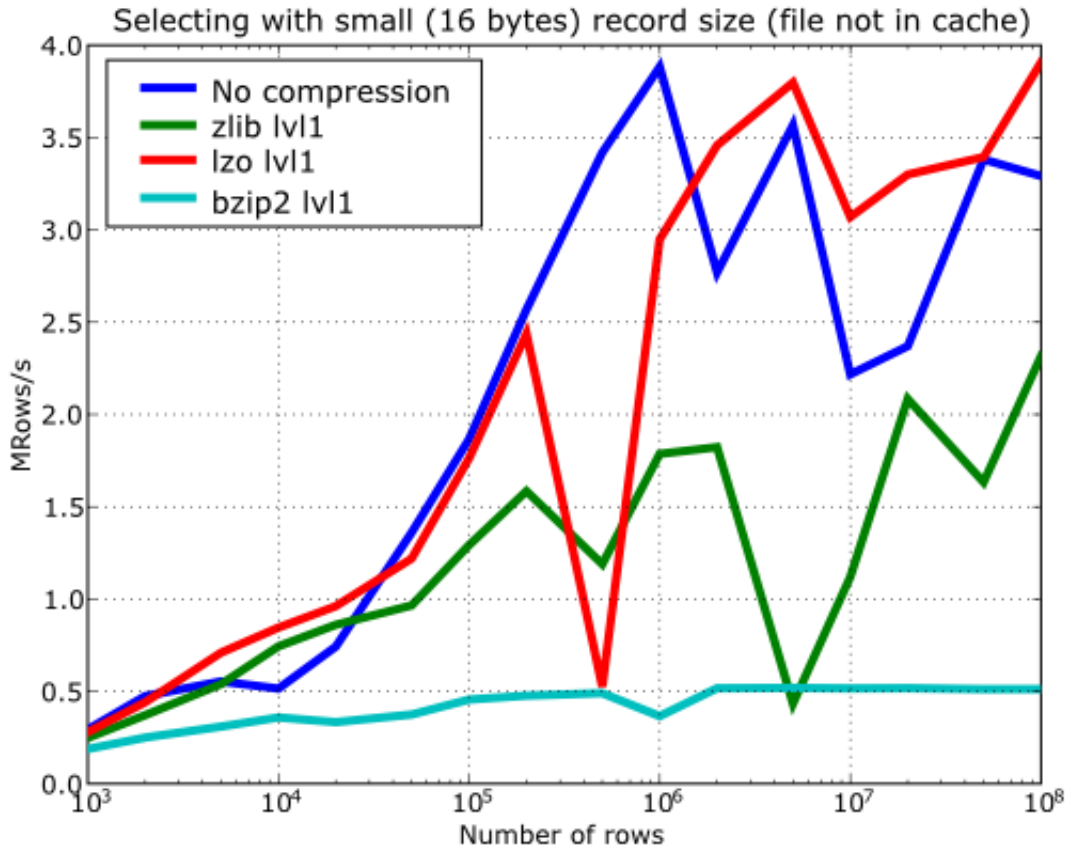


Fig. 22: Figure 16. Selecting values in tables with several compressors. The file is not in the OS cache.

So, generally speaking and looking at the experiments above, you can expect that LZO will be the fastest in both compressing and decompressing, but the one that achieves the worse compression ratio (although that may be just OK for many situations, specially when used with shuffling - see [Shuffling \(or how to make the compression process more effective\)](#)). bzip2 is the slowest, by large, in both compressing and decompressing, and besides, it does not achieve any better compression ratio than Zlib. Zlib represents a balance between them: it's somewhat slow compressing (2x) and decompressing (3x) than LZO, but it normally achieves better compression ratios.

Finally, by looking at the plots [Figure 18](#), [Figure 19](#), and the aforementioned [Figure 14](#) you can see why the recommended compression level to use for all compression libraries is 1. This is the lowest level of compression, but as the size of the underlying HDF5 chunk size is normally rather small compared with the size of compression buffers, there is not much point in increasing the latter (i.e. increasing the compression level). Nonetheless, in some situations (like for example, in extremely large tables or arrays, where the computed chunk size can be rather large) you may want to check, on your own, how the different compression levels do actually affect your application.

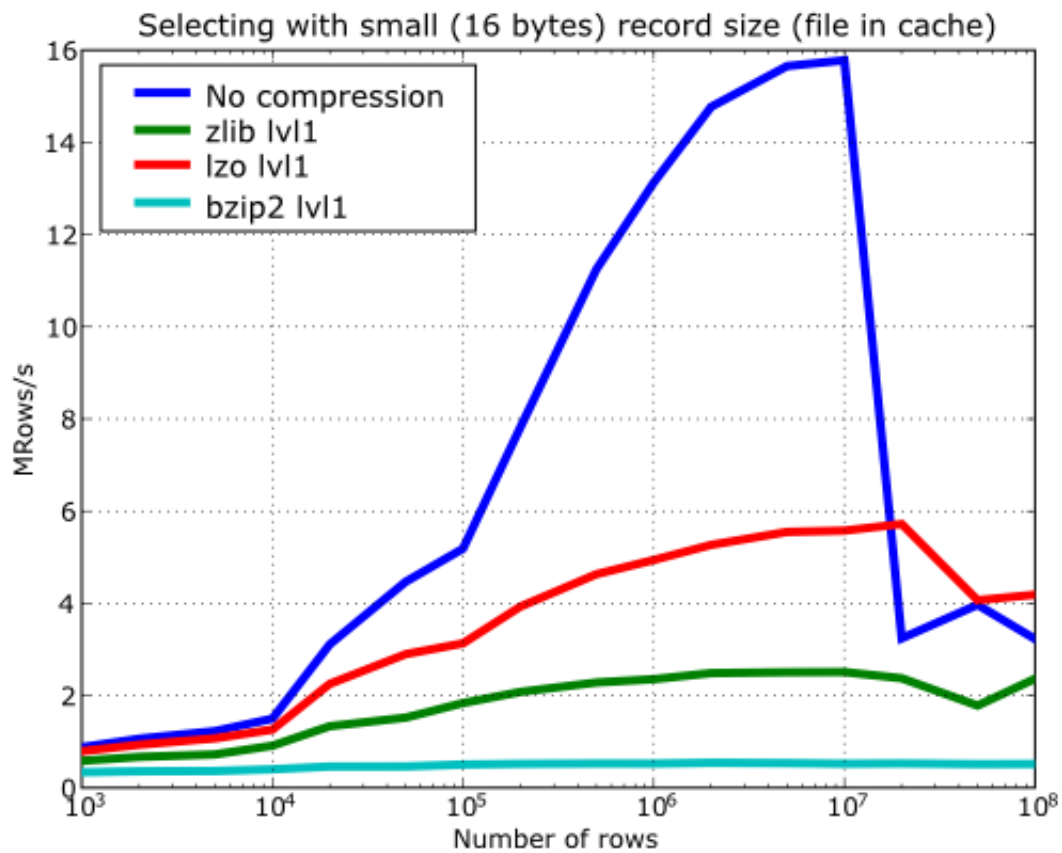


Fig. 23: Figure 17. Selecting values in tables with several compressors. The file is in the OS cache.

You can select the compression library and level by setting the `complib` and `complevel` keywords in the `Filters` class (see [The Filters class](#)). A compression level of 0 will completely disable compression (the default), 1 is the less memory and CPU time demanding level, while 9 is the maximum level and the most memory demanding and 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.

So, in conclusion, if your ultimate goal is writing and reading as fast as possible, choose LZO. If you want to reduce as much as possible your data, while retaining acceptable read speed, choose Zlib. Finally, if portability is important for you, Zlib is your best bet. So, when you want to use bzip2? Well, looking at the results, it is difficult to recommend its use in general, but you may want to experiment with it in those cases where you know that it is well suited for your data pattern (for example, for dealing with repetitive string datasets).

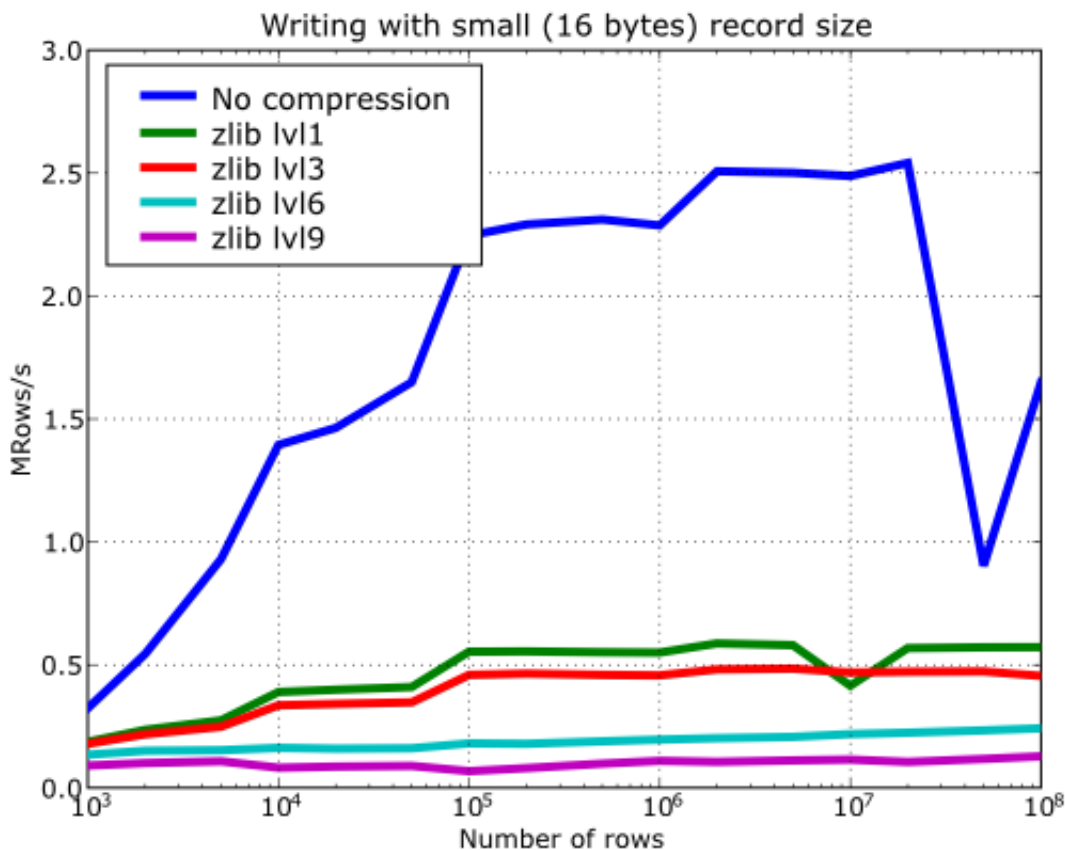


Fig. 24: Figure 18. Writing in tables with different levels of compression.

Shuffling (or how to make the compression process more effective)

The HDF5 library provides an interesting filter that can leverage the results of your favorite compressor. Its name is *shuffle*, and because it can greatly benefit compression and it does not take many CPU resources (see below for a justification), it is active *by default* in PyTables whenever compression is activated (independently of the chosen compressor). It is deactivated when compression is off (which is the default, as you already should know). Of course, you can deactivate it if you want, but this is not recommended.

Note: Since PyTables 3.3, a new *bitshuffle* filter for Blosc compressor has been added. Contrarily to *shuffle* that

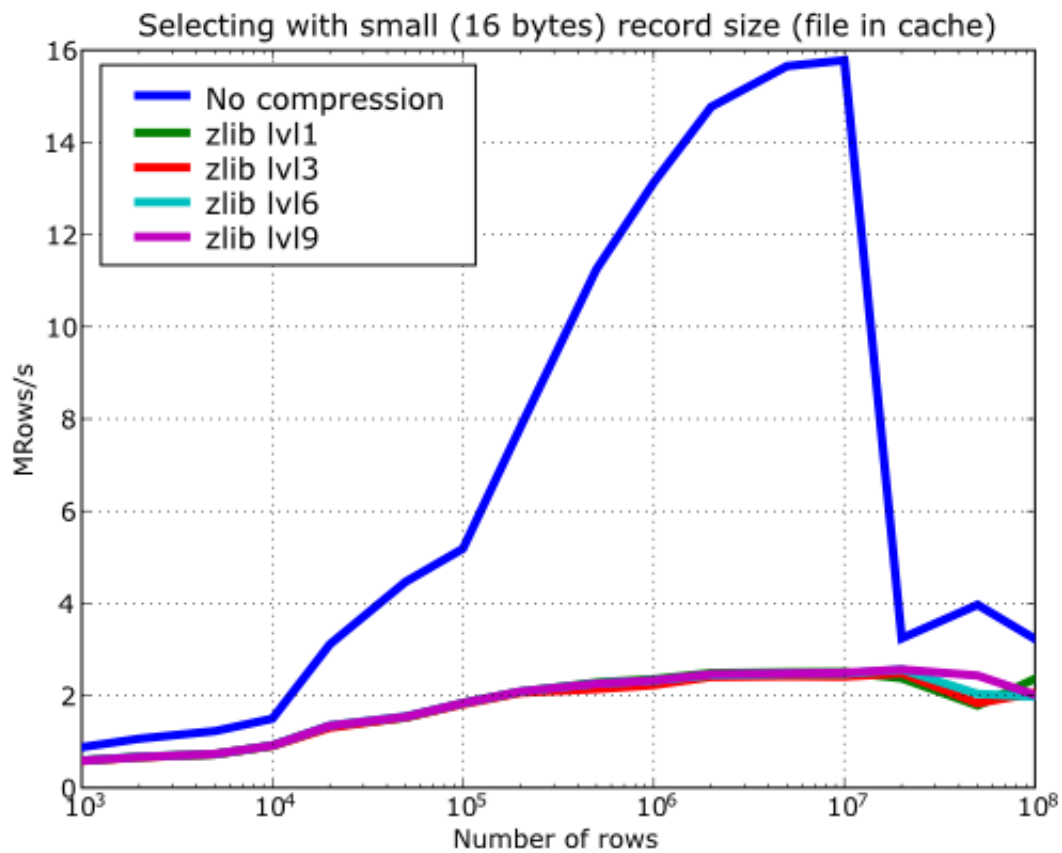


Fig. 25: Figure 19. Selecting values in tables with different levels of compression. The file is in the OS cache.

shuffles bytes, *bitshuffle* shuffles the chunk data at bit level which **could** improve compression ratios at the expense of some speed penalty. Look at the [The Filters class](#) documentation on how to activate bitshuffle and experiment with it so as to decide if it can be useful for you.

So, how does this mysterious filter exactly work? From the HDF5 reference manual:

"The shuffle filter de-interlaces a block of data by reordering the bytes. All the bytes from one consistent byte position of each data element are placed together in one block; all bytes from a second consistent byte position of each data element are placed together a second block; etc. For example, given three data elements of a 4-byte datatype stored as 012301230123, shuffling will re-order data as 000111222333. This can be a valuable step in an effective compression algorithm because the bytes in each byte position are often closely related to each other and putting them together can increase the compression ratio."

In [Figure 20](#) you can see a benchmark that shows how the *shuffle* filter can help the different libraries in compressing data. In this experiment, shuffle has made LZO compress almost 3x more (!), while Zlib and bzip2 are seeing improvements of 2x. Once again, the data for this experiment is synthetic, and *shuffle* seems to do a great work with it, but in general, the results will vary in each case³.

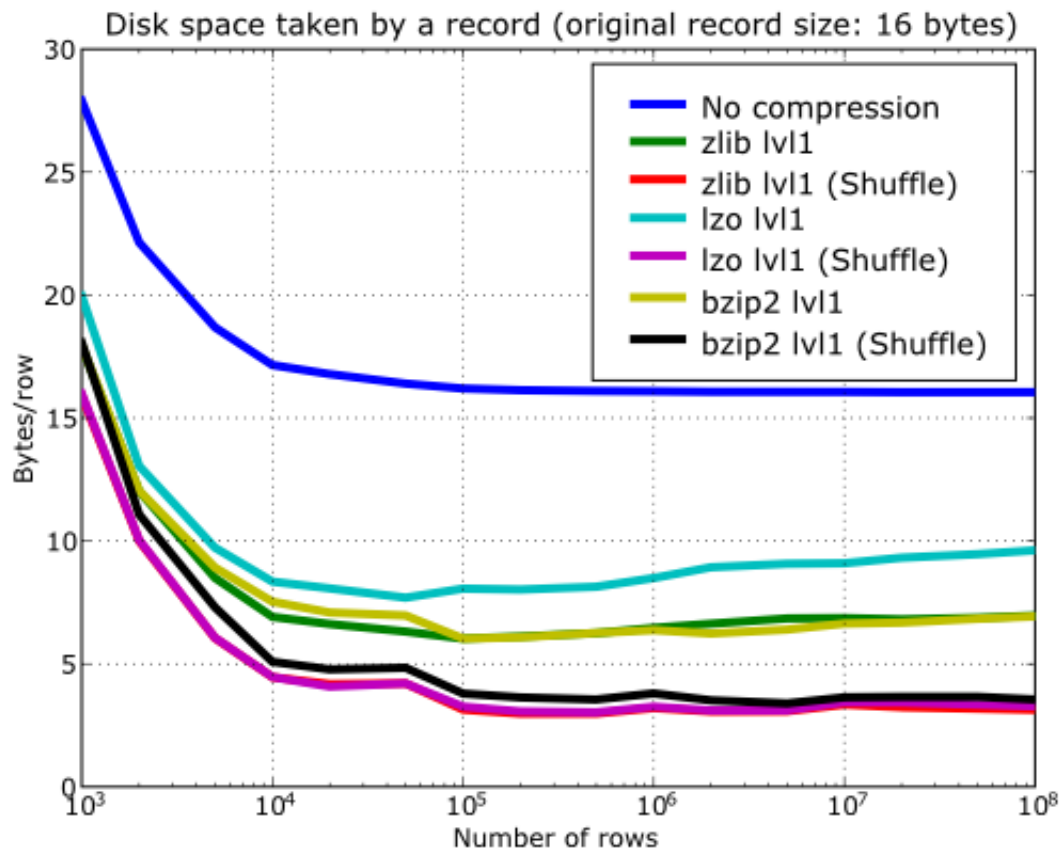


Fig. 26: **Figure 20. Comparison between different compression libraries with and without the shuffle filter.**

³ Some users reported that the typical improvement with real data is between a factor 1.5x and 2.5x over the already compressed datasets.

At any rate, the most remarkable fact about the *shuffle* filter is the relatively high level of compression that compressor filters can achieve when used in combination with it. A curious thing to note is that the Bzip2 compression rate does not seem very much improved (less than a 40%), and what is more striking, Bzip2+shuffle does compress quite *less* than Zlib+shuffle or LZO+shuffle combinations, which is kind of unexpected. The thing that seems clear is that Bzip2 is not very good at compressing patterns that result of shuffle application. As always, you may want to experiment with your own data before widely applying the Bzip2+shuffle combination in order to avoid surprises.

Now, how does shuffling affect performance? Well, if you look at plots [Figure 21](#), [Figure 22](#) and [Figure 23](#), you will get a somewhat unexpected (but pleasant) surprise. Roughly, *shuffle* makes the writing process (shuffling+compressing) faster (approximately a 15% for LZO, 30% for Bzip2 and a 80% for Zlib), which is an interesting result by itself. But perhaps more exciting is the fact that the reading process (unshuffling+decompressing) is also accelerated by a similar extent (a 20% for LZO, 60% for Zlib and a 75% for Bzip2, roughly).

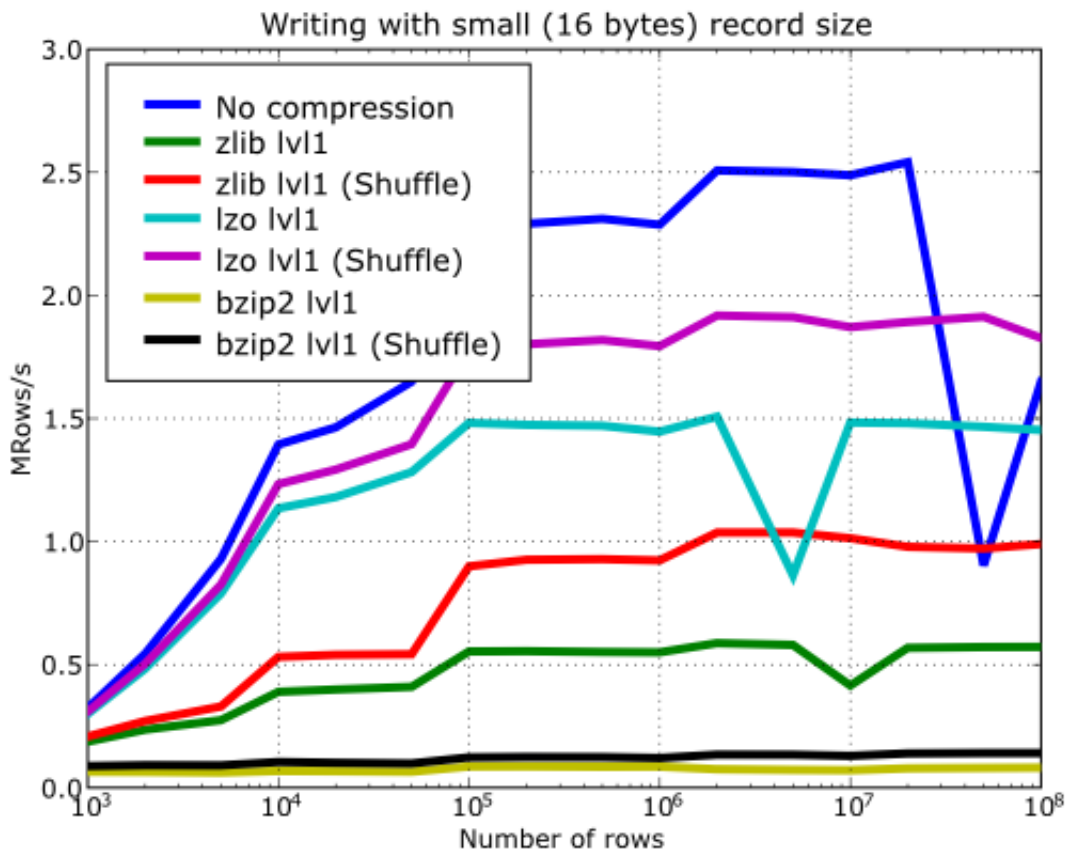


Fig. 27: **Figure 21. Writing with different compression libraries with and without the shuffle filter.**

You may wonder why introducing another filter in the write/read pipelines does effectively accelerate the throughput. Well, maybe data elements are more similar or related column-wise than row-wise, i.e. contiguous elements in the same column are more alike, so shuffling makes the job of the compressor easier (faster) and more effective (greater ratios). As a side effect, compressed chunks do fit better in the CPU cache (at least, the chunks are smaller!) so that the process of unshuffle/decompress can make a better use of the cache (i.e. reducing the number of CPU cache faults).

So, given the potential gains (faster writing and reading, but specially much improved compression level), it is a good thing to have such a filter enabled by default in the battle for discovering redundancy when you want to compress your data, just as PyTables does.

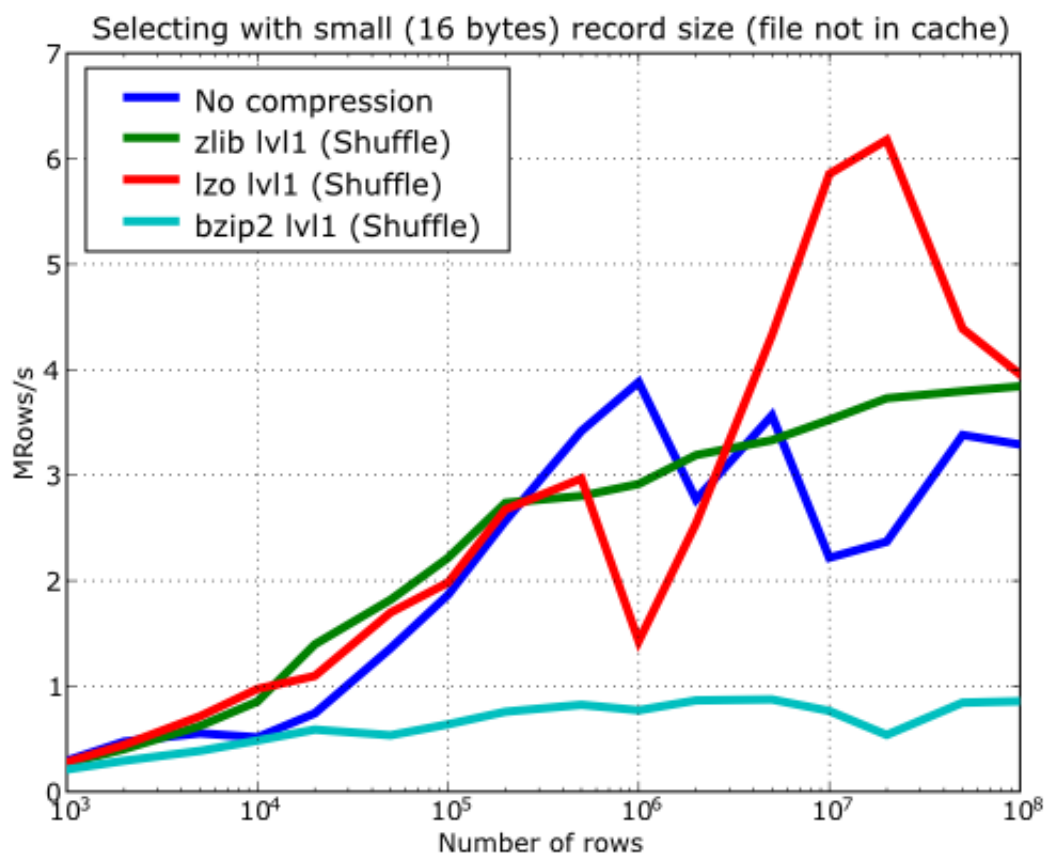


Fig. 28: **Figure 22.** Reading with different compression libraries with the shuffle filter. The file is not in OS cache.

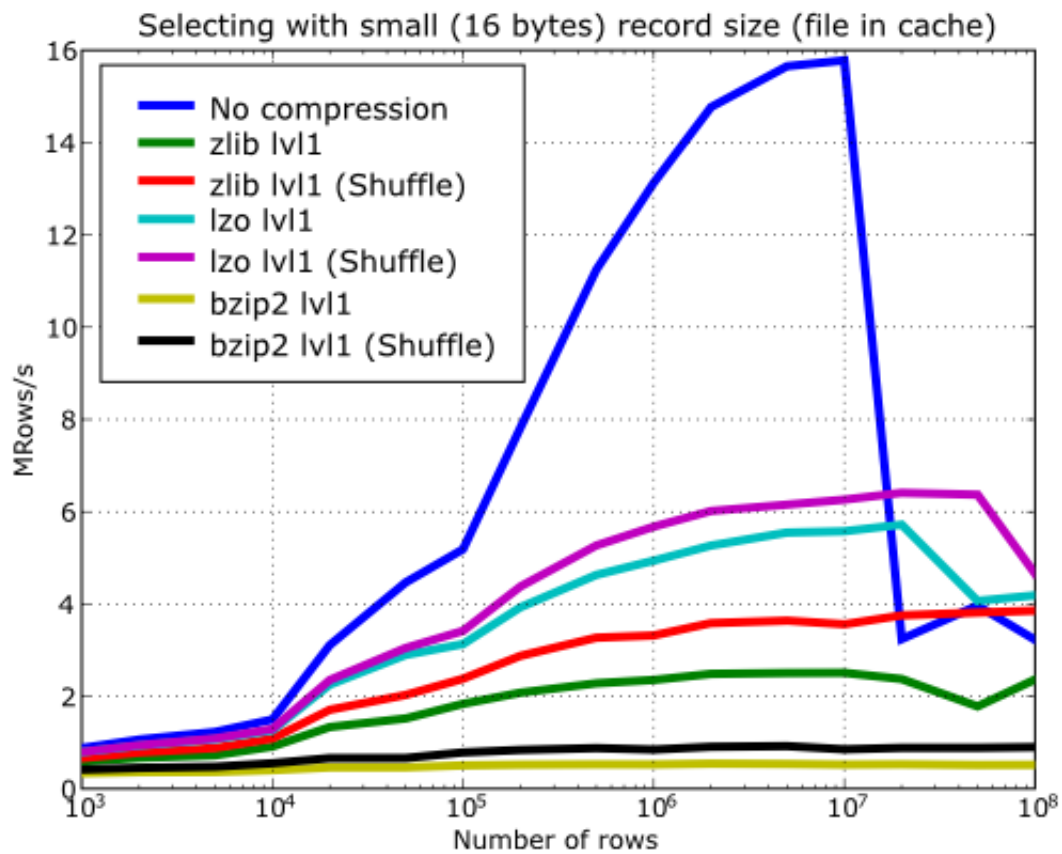


Fig. 29: **Figure 23. Reading with different compression libraries with and without the shuffle filter. The file is in OS cache.**

1.5.4 Using Psyco

Psyco (see [\[PSYCO\]](#)) 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 emits 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 give it a try. However, it only runs on Intel 386 architectures, so if you are using other architectures, you are out of luck (and, moreover, it seems that there are no plans to support other platforms). Besides, with the addition of flexible (and very fast) in-kernel queries (by the way, they cannot be optimized at all by Psyco), the use of Psyco will only help in rather few scenarios. In fact, the only important situation that you might benefit right now from using Psyco (I mean, in PyTables contexts) is for speeding-up the write speed in tables when using the Row interface (see [The Row class](#)). But again, this latter case can also be accelerated by using the `Table.append()` method and building your own buffers⁴.

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

```
def read_file(filename):
    "Select data from all the tables in filename"
    fileh = open_file(filename, mode = "r")
    result = []
    for table in fileh("/", 'Table'):
        result = [p['var3'] for p in table if p['var2'] <= 20]
    fileh.close()
    return result

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

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

```
if __name__=="__main__":
    import psyco
    psyco.bind(read_file)
    print(read_file("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.

You can see in the graphs [Figure 24](#) and [Figure 25](#) 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 (from a thousand to a million rows), 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 [Compression issues](#)). As always, there is no substitute for experimentation with your own dataset.

1.5.5 Getting the most from the node LRU cache

One limitation of the initial versions of PyTables was that they needed to load all nodes in a file completely before being ready to deal with them, making the opening times for files with a lot of nodes very high and unacceptable in many cases.

Starting from PyTables 1.2 on, a new lazy node loading schema was setup that avoids loading all the nodes of the *object tree* in memory. In addition, a new LRU cache was introduced in order to accelerate the access to already visited nodes. This cache (one per file) is responsible for keeping up the most recently visited nodes in memory and

⁴ So, there is not much point in using Psyco with recent versions of PyTables anymore.

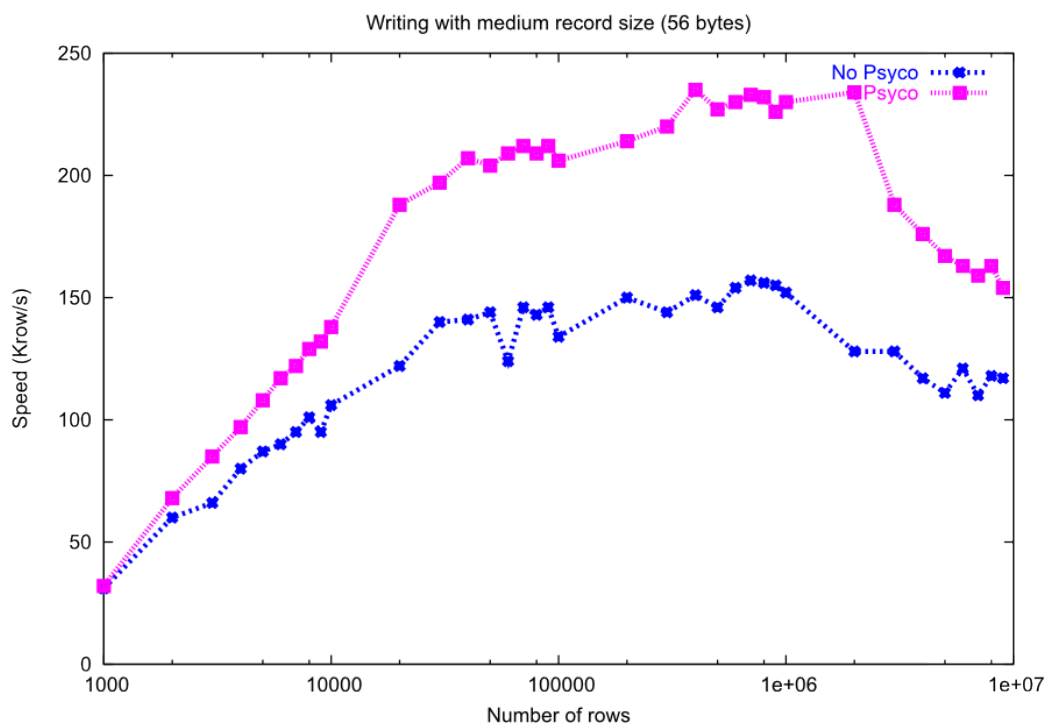


Fig. 30: **Figure 24. Writing tables with/without Psyco.**

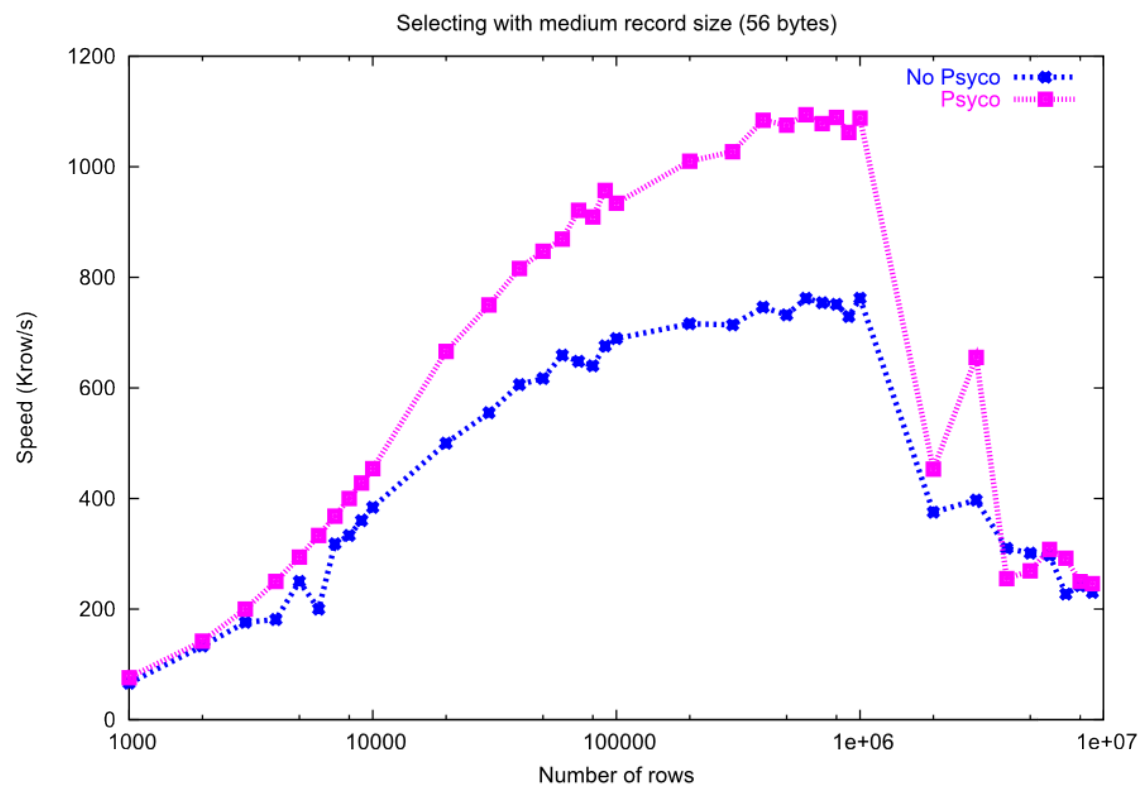


Fig. 31: Figure 25. Reading tables with/without Psycopy.

discard the least recent used ones. This represents a big advantage over the old schema, not only in terms of memory usage (as there is no need to load *every* node in memory), but it also adds very convenient optimizations for working interactively like, for example, speeding-up the opening times of files with lots of nodes, allowing to open almost any kind of file in typically less than one tenth of second (compare this with the more than 10 seconds for files with more than 10000 nodes in PyTables pre-1.2 era) as well as optimizing the access to frequently visited nodes. See for more info on the advantages (and also drawbacks) of this approach.

One thing that deserves some discussion is the election of the parameter that sets the maximum amount of nodes to be kept in memory at any time. As PyTables is meant to be deployed in machines that can have potentially low memory, the default for it is quite conservative (you can look at its actual value in the `parameters.NODE_CACHE_SLOTS` parameter in module `tables/parameters.py`). However, if you usually need to deal with files that have many more nodes than the maximum default, and you have a lot of free memory in your system, then you may want to experiment in order to see which is the appropriate value of `parameters.NODE_CACHE_SLOTS` that fits better your needs.

As an example, look at the next code:

```
def browse_tables(filename):
    fileh = open_file(filename, 'a')
    group = fileh.root.newgroup
    for j in range(10):
        for tt in fileh.walk_nodes(group, "Table"):
            title = tt.attrs.TITLE
            for row in tt:
                pass
    fileh.close()
```

We will be running the code above against a couple of files having a `/newgroup` containing 100 tables and 1000 tables respectively. In addition, this benchmark is run twice for two different values of the LRU cache size, specifically 256 and 1024. You can see the results in [table](#).

Table 1: Retrieval speed and memory consumption depending on the number of nodes in LRU cache.

Number:		100 nodes				1000 nodes			
Mem and Speed		Memory (MB)		Time (ms)		Memory (MB)		Time (ms)	
Node is coming from...	Cache size	256	1024	256	1024	256	1024	256	1024
Disk		14	14	1.24	1.24	51	66	1.33	1.31
Cache		14	14	0.53	0.52	65	73	1.35	0.68

From the data in [table](#), one can see that when the number of objects that you are dealing with does fit in cache, you will get better access times to them. Also, incrementing the node cache size effectively consumes more memory *only* if the total nodes exceeds the slots in cache; otherwise the memory consumption remains the same. It is also worth noting that incrementing the node cache size in the case you want to fit all your nodes in cache does not take much more memory than being too conservative. On the other hand, it might happen that the speed-up that you can achieve by allocating more slots in your cache is not worth the amount of memory used.

Also worth noting is that if you have a lot of memory available and performance is absolutely critical, you may want to try out a negative value for `parameters.NODE_CACHE_SLOTS`. This will cause that all the touched nodes will be kept in an internal dictionary and this is the faster way to load/retrieve nodes. However, and in order to avoid a large memory consumption, the user will be warned when the number of loaded nodes will reach the `-NODE_CACHE_SLOTS` value.

Finally, a value of zero in `parameters.NODE_CACHE_SLOTS` means that any cache mechanism is disabled.

At any rate, if you feel that this issue is important for you, there is no replacement for setting your own experiments up in order to proceed to fine-tune the `parameters.NODE_CACHE_SLOTS` parameter.

Note: PyTables >= 2.3 sports an optimized LRU cache node written in C, so you should expect significantly faster

LRU cache operations when working with it.

Note: Numerical results reported in *table* have been obtained with PyTables < 3.1. In PyTables 3.1 the node cache mechanism has been completely redesigned so while all comments above are still valid, numerical values could be a little bit different from the ones reported in *table*.

1.5.6 Compacting your PyTables files

Let's suppose that you have a file where you have made a lot of row deletions on one or more tables, or deleted many leaves or even entire subtrees. These operations might leave *holes* (i.e. space that is not used anymore) in your files that may potentially affect not only the size of the files but, more importantly, the performance of I/O. This is because when you delete a lot of rows in a table, the space is not automatically recovered on the fly. In addition, if you add many more rows to a table than specified in the `expectedrows` keyword at creation time this may affect performance as well, as explained in *Informing PyTables about expected number of rows in tables or arrays*.

In order to cope with these issues, you should be aware that PyTables includes a handy utility called `ptrepack` which can be very useful not only to compact *fragmented* files, but also to adjust some internal parameters in order to use better buffer and chunk sizes for optimum I/O speed. Please check the *ptrepack* for a brief tutorial on its use.

Another thing that you might want to use `ptrepack` for is changing the compression filters or compression levels on your existing data for different goals, like checking how this can affect both final size and I/O performance, or getting rid of the optional compressors like LZO or bzip2 in your existing files, in case you want to use them with generic HDF5 tools that do not have support for these filters.

COMPLEMENTARY MODULES

2.1 filenode - simulating a filesystem with PyTables

2.1.1 What is filenode?

filenode is a module which enables you to create a PyTables database of nodes which can be used like regular opened files in Python. In other words, you can store a file in a PyTables database, and read and write it as you would do with any other file in Python. Used in conjunction with PyTables hierarchical database organization, you can have your database turned into an open, extensible, efficient, high capacity, portable and metadata-rich filesystem for data exchange with other systems (including backup purposes).

Between the main features of filenode, one can list:

- *Open:* Since it relies on PyTables, which in turn, sits over HDF5 (see [\[HDGG1\]](#)), a standard hierarchical data format from NCSA.
- *Extensible:* You can define new types of nodes, and their instances will be safely preserved (as are normal groups, leafs and attributes) by PyTables applications having no knowledge of their types. Moreover, the set of possible attributes for a node is not fixed, so you can define your own node attributes.
- *Efficient:* Thanks to PyTables' proven extreme efficiency on handling huge amounts of data, filenode can make use of PyTables' on-the-fly compression and decompression of data.
- *High capacity:* Since PyTables and HDF5 are designed for massive data storage (they use 64-bit addressing even where the platform does not support it natively).
- *Portable:* Since the HDF5 format has an architecture-neutral design, and the HDF5 libraries and PyTables are known to run under a variety of platforms. Besides that, a PyTables database fits into a single file, which poses no trouble for transportation.
- *Metadata-rich:* Since PyTables can store arbitrary key-value pairs (even Python objects!) for every database node. Metadata may include authorship, keywords, MIME types and encodings, ownership information, access control lists (ACL), decoding functions and anything you can imagine!

2.1.2 Finding a filenode node

filenode nodes can be recognized because they have a `NODE_TYPE` system attribute with a 'file' value. It is recommended that you use the `File.get_node_attr()` method of `tables.File` class to get the `NODE_TYPE` attribute independently of the nature (group or leaf) of the node, so you do not need to care about.

2.1.3 filenode - simulating files inside PyTables

The filenode module is part of the nodes sub-package of PyTables. The recommended way to import the module is:

```
>>> from tables.nodes import filenode
```

However, `filenode` exports very few symbols, so you can import `*` for interactive usage. In fact, you will most probably only use the `NodeType` constant and the `new_node()` and `open_node()` calls.

The `NodeType` constant contains the value that the `NODE_TYPE` system attribute of a node file is expected to contain ('file', as we have seen). Although this is not expected to change, you should use `filenode.NodeType` instead of the literal 'file' when possible.

`new_node()` and `open_node()` are the equivalent to the Python `file()` call (alias `open()`) for ordinary files. Their arguments differ from that of `file()`, but this is the only point where you will note the difference between working with a node file and working with an ordinary file.

For this little tutorial, we will assume that we have a PyTables database opened for writing. Also, if you are somewhat lazy at typing sentences, the code that we are going to explain is included in the `examples/filenodes1.py` file.

You can create a brand new file with these sentences:

```
>>> import tables
>>> h5file = tables.open_file('fnode.h5', 'w')
```

Creating a new file node

Creation of a new file node is achieved with the `new_node()` call. You must tell it in which PyTables file you want to create it, where in the PyTables hierarchy you want to create the node and which will be its name. The PyTables file is the first argument to `new_node()`; it will be also called the 'host PyTables file'. The other two arguments must be given as keyword arguments `where` and `name`, respectively. As a result of the call, a brand new appendable and readable file node object is returned.

So let us create a new node file in the previously opened `h5file` PyTables file, named 'fnode_test' and placed right under the root of the database hierarchy. This is that command:

```
>>> fnode = filenode.new_node(h5file, where='/', name='fnode_test')
```

That is basically all you need to create a file node. Simple, isn't it? From that point on, you can use `fnode` as any opened Python file (i.e. you can write data, read data, lines of text and so on).

`new_node()` accepts some more keyword arguments. You can give a title to your file with the `title` argument. You can use PyTables' compression features with the `filters` argument. If you know beforehand the size that your file will have, you can give its final file size in bytes to the `expectedsize` argument so that the PyTables library would be able to optimize the data access.

`new_node()` creates a PyTables node where it is told to. To prove it, we will try to get the `NODE_TYPE` attribute from the newly created node:

```
>>> print(h5file.get_node_attr('/fnode_test', 'NODE_TYPE'))
file
```

Using a file node

As stated above, you can use the new node file as any other opened file. Let us try to write some text in and read it:

```
>>> print("This is a test text line.", file=fnode)
>>> print("And this is another one.", file=fnode)
>>> print(file=fnode)
```

(continues on next page)

(continued from previous page)

```
>>> fnode.write("Of course, file methods can also be used.")
>>>
>>> fnode.seek(0)  # Go back to the beginning of file.
>>>
>>> for line in fnode:
...     print(repr(line))
'This is a test text line.\n'
'And this is another one.\n'
'\n'
'Of course, file methods can also be used.'
```

This was run on a Unix system, so newlines are expressed as ‘\n’. In fact, you can override the line separator for a file by setting its `line_separator` property to any string you want.

While using a file node, you should take care of closing it *before* you close the PyTables host file. Because of the way PyTables works, your data it will not be at a risk, but every operation you execute after closing the host file will fail with a `ValueError`. To close a file node, simply delete it or call its `close()` method:

```
>>> fnode.close()
>>> print(fnode.closed)
True
```

Opening an existing file node

If you have a file node that you created using `new_node()`, you can open it later by calling `open_node()`. Its arguments are similar to that of `file()` or `open()`: the first argument is the PyTables node that you want to open (i.e. a node with a `NODE_TYPE` attribute having a ‘file’ value), and the second argument is a mode string indicating how to open the file. Contrary to `file()`, `open_node()` can not be used to create a new file node.

File nodes can be opened in read-only mode (‘r’) or in read-and-append mode (‘a+’). Reading from a file node is allowed in both modes, but appending is only allowed in the second one. Just like Python files do, writing data to an appendable file places it after the file pointer if it is on or beyond the end of the file, or otherwise after the existing data. Let us see an example:

```
>>> node = h5file.root.fnode_test
>>> fnode = filenode.open_node(node, 'a+')
>>> print(repr(fnode.readline()))
'This is a test text line.\n'
>>> print(fnode.tell())
26
>>> print("This is a new line.", file=fnode)
>>> print(repr(fnode.readline()))
''
```

Of course, the data append process places the pointer at the end of the file, so the last `readline()` call hit EOF. Let us seek to the beginning of the file to see the whole contents of our file:

```
>>> fnode.seek(0)
>>> for line in fnode:
...     print(repr(line))
'This is a test text line.\n'
'And this is another one.\n'
'\n'
'Of course, file methods can also be used.This is a new line.\n'
```

As you can check, the last string we wrote was correctly appended at the end of the file, instead of overwriting the second line, where the file pointer was positioned by the time of the appending.

Adding metadata to a file node

You can associate arbitrary metadata to any open node file, regardless of its mode, as long as the host PyTables file is writable. Of course, you could use the `set_node_attr()` method of `tables.File` to do it directly on the proper node, but `filenode` offers a much more comfortable way to do it. `filenode` objects have an `attrs` property which gives you direct access to their corresponding `AttributeSet` object.

For instance, let us see how to associate MIME type metadata to our file node:

```
>>> fnode.attrs.content_type = 'text/plain; charset=us-ascii'
```

As simple as A-B-C. You can put nearly anything in an attribute, which opens the way to authorship, keywords, permissions and more. Moreover, there is not a fixed list of attributes. However, you should avoid names in all caps or starting with `'_'`, since PyTables and `filenode` may use them internally. Some valid examples:

```
>>> fnode.attrs.author = "Ivan Vilata i Balaguer"
>>> fnode.attrs.creation_date = '2004-10-20T13:25:25+0200'
>>> fnode.attrs.keywords_en = ["FileName", "test", "metadata"]
>>> fnode.attrs.keywords_ca = ["FileName", "prova", "metadades"]
>>> fnode.attrs.owner = 'ivan'
>>> fnode.attrs.acl = {'ivan': 'rw', '@users': 'r'}
```

You can check that these attributes get stored by running the `ptdump` command on the host PyTables file.

```
$ ptdump -a fnode.h5:/fnode_test
/fnode_test (EArray(113,)) ''
/fnode_test.attrs (AttributeSet), 14 attributes:
[CLASS := 'EARRAY',
EXTDIM := 0,
FLAVOR := 'numpy',
NODE_TYPE := 'file',
NODE_TYPE_VERSION := 2,
TITLE := '',
VERSION := '1.2',
acl := {'ivan': 'rw', '@users': 'r'},
author := 'Ivan Vilata i Balaguer',
content_type := 'text/plain; charset=us-ascii',
creation_date := '2004-10-20T13:25:25+0200',
keywords_ca := ['FileName', 'prova', 'metadades'],
keywords_en := ['FileName', 'test', 'metadata'],
owner := 'ivan']
```

Note that `filenode` makes no assumptions about the meaning of your metadata, so its handling is entirely left to your needs and imagination.

2.1.4 Complementary notes

You can use file nodes and PyTables groups to mimic a filesystem with files and directories. Since you can store nearly anything you want as file metadata, this enables you to use a PyTables file as a portable compressed backup, even between radically different platforms. Take this with a grain of salt, since node files are restricted in their naming (only valid Python identifiers are valid); however, remember that you can use node titles and metadata to overcome this limitation. Also, you may need to devise some strategy to represent special files such as devices, sockets and such (not necessarily using `filenode`).

We are eager to hear your opinion about filenode and its potential uses. Suggestions to improve filenode and create other node types are also welcome. Do not hesitate to contact us!

2.1.5 Current limitations

filenode is still a young piece of software, so it lacks some functionality. This is a list of known current limitations:

1. Node files can only be opened for read-only or read and append mode. This should be enhanced in the future.
2. Near future?
3. Only binary I/O is supported currently (read/write strings of bytes)
4. There is no universal newline support yet. The only new-line character used at the moment is `\n`. This is likely to be improved in a near future.
5. Sparse files (files with lots of zeros) are not treated specially; if you want them to take less space, you should be better off using compression.

These limitations still make filenode entirely adequate to work with most binary and text files. Of course, suggestions and patches are welcome.

See *Filenode Module* for detailed documentation on the filenode interface.

3.1 Supported data types in PyTables

All PyTables datasets can handle the complete set of data types supported by the NumPy (see [\[NUMPY\]](#)) package in Python. The data types for table fields can be set via instances of the `Col` class and its descendants (see [The `Col` class and its descendants](#)), while the data type of array elements can be set through the use of the `Atom` class and its descendants (see [The `Atom` class and its descendants](#)).

PyTables uses ordinary strings to represent its *types*, with most of them matching the names of NumPy scalar types. Usually, a PyTables type consists of two parts: a *kind* and a *precision* in bits. The precision may be omitted in types with just one supported precision (like `bool`) or with a non-fixed size (like `string`).

There are eight kinds of types supported by PyTables:

- `bool`: Boolean (`true/false`) types. Supported precisions: 8 (default) bits.
- `int`: Signed integer types. Supported precisions: 8, 16, 32 (default) and 64 bits.
- `uint`: Unsigned integer types. Supported precisions: 8, 16, 32 (default) and 64 bits.
- `float`: Floating point types. Supported precisions: 16, 32, 64 (default) bits and extended precision floating point (see [note on floating point types](#)).
- `complex`: Complex number types. Supported precisions: 64 (32+32), 128 (64+64, default) bits and extended precision complex (see [note on floating point types](#)).
- `string`: Raw string types. Supported precisions: 8-bit positive multiples.
- `time`: Data/time types. Supported precisions: 32 and 64 (default) bits.
- `enum`: Enumerated types. Precision depends on base type.

Note: Floating point types.

The half precision floating point data type (`float16`) and extended precision ones (`float96`, `float128`, `complex192`, `complex256`) are only available if `numpy` supports them on the host platform.

Also, in order to use the half precision floating point type (`float16`) it is required `numpy >= 1.6.0`.

The `time` and `enum` kinds are a little bit special, since they represent HDF5 types which have no direct Python counterpart, though atoms of these kinds have a more-or-less equivalent NumPy data type.

There are two types of `time`: 4-byte signed integer (`time32`) and 8-byte double precision floating point (`time64`). Both of them reflect the number of seconds since the Unix epoch, i.e. Jan 1 00:00:00 UTC 1970. They are stored in memory as NumPy's `int32` and `float64`, respectively, and in the HDF5 file using the `H5T_TIME` class. Integer times are stored on disk as such, while floating point times are split into two signed integer values representing seconds and microseconds (beware: smaller decimals will be lost!).

PyTables also supports HDF5 H5T_ENUM *enumerations* (restricted sets of unique name and unique value pairs). The NumPy representation of an enumerated value (an Enum, see [The Enum class](#)) depends on the concrete *base type* used to store the enumeration in the HDF5 file. Currently, only scalar integer values (both signed and unsigned) are supported in enumerations. This restriction may be lifted when HDF5 supports other kinds of enumerated values.

Here you have a quick reference to the complete set of supported data types:

Table 1: **Data types supported for array elements and tables columns in PyTables.**

Type Code	Description	C Type	Size (in bytes)	Python Counter-part
bool	boolean	unsigned char	1	bool
int8	8-bit integer	signed char	1	int
uint8	8-bit unsigned integer	unsigned char	1	int
int16	16-bit integer	short	2	int
uint16	16-bit unsigned integer	unsigned short	2	int
int32	integer	int	4	int
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
float16 ¹	half-precision float	.	2	.
float32	single-precision float	float	4	float
float64	double-precision float	double	8	float
float96 ¹²	extended precision float	.	12	.
float128 ¹²	extended precision float	.	16	.
complex64	single-precision complex	struct { float r, i; }	8	complex
complex128	double-precision complex	struct { double r, i; }	16	complex
complex192 ¹	extended precision complex	.	24	.
complex256 ¹	extended precision complex	.	32	.
string	arbitrary length string	char[]	.	str
time32	integer time	POSIX's time_t	4	int
time64	floating point time	POSIX's struct timeval	8	float
enum	enumerated value	enum	.	.

¹ see the above [note on floating point types](#).

² currently in `numpy`. “float96” and “float128” are equivalent of “longdouble” i.e. 80 bit extended precision floating point.

3.2 Condition Syntax

Conditions in PyTables are used in methods related with in-kernel and indexed searches such as `Table.where()` or `Table.read_where()`. They are interpreted using Numexpr, a powerful package for achieving C-speed computation of array operations (see [\[NUMEXPR\]](#)).

A condition on a table is just a *string* containing a Python expression involving *at least one column*, and maybe some constants and external variables, all combined with algebraic operators and functions. The result of a valid condition is always a *boolean array* of the same length as the table, where the *i*-th element is true if the value of the expression on the *i*-th row of the table evaluates to true.

That is the reason why multidimensional fields in a table are not supported in conditions, since the truth value of each resulting multidimensional boolean value is not obvious. Usually, a method using a condition will only consider the rows where the boolean result is true.

For instance, the condition `'sqrt(x*x + y*y) < 1'` applied on a table with `x` and `y` columns consisting of floating point numbers results in a boolean array where the *i*-th element is true if (unsurprisingly) the value of the square root of the sum of squares of `x` and `y` is less than 1. The `sqrt()` function works element-wise, the `1` constant is adequately broadcast to an array of ones of the length of the table for evaluation, and the *less than* operator makes the result a valid boolean array. A condition like `'mycolumn'` alone will not usually be valid, unless `mycolumn` is itself a column of scalar, boolean values.

In the previous conditions, `mycolumn`, `x` and `y` are examples of *variables* which are associated with columns. Methods supporting conditions do usually provide their own ways of binding variable names to columns and other values. You can read the documentation of `Table.where()` for more information on that. Also, please note that the names `None`, `True` and `False`, besides the names of functions (see below) *can not be overridden*, but you can always define other new names for the objects you intend to use.

Values in a condition may have the following types:

- 8-bit boolean (`bool`).
- 32-bit signed integer (`int`).
- 64-bit signed integer (`long`).
- 32-bit, single-precision floating point number (`float` or `float32`).
- 64-bit, double-precision floating point number (`double` or `float64`).
- 2x64-bit, double-precision complex number (`complex`).
- Raw string of bytes (`str`).

Nevertheless, if the type passed is not among the above ones, it will be silently upcasted, so you don't need to worry too much about passing supported types, except for the Unsigned 64 bits integer, that cannot be upcasted to any of the supported types.

However, the types in PyTables conditions are somewhat stricter than those of Python. For instance, the *only* valid constants for booleans are `True` and `False`, and they are *never* automatically cast to integers. The type strengthening also affects the availability of operators and functions. Beyond that, the usual type inference rules apply.

Conditions support the set of operators listed below:

- Logical operators: `&`, `|`, `~`.
- Comparison operators: `<`, `<=`, `==`, `!=`, `>=`, `>`.
- Unary arithmetic operators: `-`.
- Binary arithmetic operators: `+`, `-`, `*`, `/`, `**`, `%`.

Types do not support all operators. Boolean values only support logical and strict (in)equality comparison operators, while strings only support comparisons, numbers do not work with logical operators, and complex comparisons can only check for strict (in)equality. Unsupported operations (including invalid castings) raise `NotImplementedError` exceptions.

You may have noticed the special meaning of the usually bitwise operators `&`, `|` and `~`. Because of the way Python handles the short-circuiting of logical operators and the truth values of their operands, conditions must use the bitwise operator equivalents instead. This is not difficult to remember, but you must be careful because bitwise operators have a *higher precedence* than logical operators. For instance, `'a and b == c'` (*a is true AND b is equal to c*) is *not* equivalent to `'a & b == c'` (*a AND b is equal to c*). The safest way to avoid confusions is to use *parentheses* around logical operators, like this: `'a & (b == c)'`. Another effect of short-circuiting is that expressions like `'0 < x < 1'` will *not* work as expected; you should use `'(0 < x) & (x < 1)'`.

All of this may be solved if Python supported overloadable boolean operators (see PEP 335) or some kind of non-shortcircuiting boolean operators (like C's `&&`, `||` and `!`).

You can also use the following functions in conditions:

- `where(bool, number1, number2)`: number - number1 if the bool condition is true, number2 otherwise.
- `{sin,cos,tan}(float|complex)`: float|complex - trigonometric sine, cosine or tangent.
- `{arcsin,arccos,arctan}(float|complex)`: float|complex - trigonometric inverse sine, cosine or tangent.
- `arctan2(float1, float2)`: float - trigonometric inverse tangent of float1/float2.
- `{sinh,cosh,tanh}(float|complex)`: float|complex - hyperbolic sine, cosine or tangent.
- `{arcsinh,arccosh,arctanh}(float|complex)`: float|complex - hyperbolic inverse sine, cosine or tangent.
- `{log,log10,log1p}(float|complex)`: float|complex - natural, base-10 and $\log(1+x)$ logarithms.
- `{exp,expm1}(float|complex)`: float|complex - exponential and exponential minus one.
- `sqrt(float|complex)`: float|complex - square root.
- `abs(float|complex)`: float|complex - absolute value.
- `{real,imag}(complex)`: float - real or imaginary part of complex.
- `complex(float, float)`: complex - complex from real and imaginary parts.

3.3 PyTables parameter files

PyTables issues warnings when certain limits are exceeded. Those limits are not intrinsic limitations of the underlying software, but rather are proactive measures to avoid large resource consumptions. The default limits should be enough for most of cases, and users should try to respect them. However, in some situations, it can be convenient to increase (or decrease) these limits.

Also, and in order to get maximum performance, PyTables implements a series of sophisticated features, like I/O buffers or different kind of caches (for nodes, chunks and other internal metadata). These features comes with a default set of parameters that ensures a decent performance in most of situations. But, as there is always a need for every case, it is handy to have the possibility to fine-tune some of these parameters.

Because of these reasons, PyTables implements a couple of ways to change the values of these parameters. All the *tunable* parameters live in the `tables/parameters.py`. The user can choose to change them in the parameter files themselves for a global and persistent change. Moreover, if he wants a finer control, he can pass any of these parameters directly to the `tables.open_file()` function, and the new parameters will only take effect in the corresponding file (the defaults will continue to be in the parameter files).

A description of all of the tunable parameters follows. As the defaults stated here may change from release to release, please check with your actual parameter files so as to know your actual default values.

Warning: Changing the next parameters may have a very bad effect in the resource consumption and performance of your PyTables scripts.

Please be careful when touching these!

3.3.1 Tunable parameters in parameters.py

Recommended maximum values

Cache limits

Parameters for the different internal caches

Parameters for general cache behaviour

Warning: The next parameters will not take any effect if passed to the `open_file()` function, so they can only be changed in a *global* way. You can change them in the file, but this is strongly discouraged unless you know well what you are doing.

Parameters for the I/O buffer in Leaf objects

Miscellaneous

HDF5 driver management

3.4 Utilities

PyTables comes with a couple of utilities that make the life easier to the user. One is called `ptdump` and lets you see the contents of a PyTables file (or generic HDF5 file, if supported). The other one is named `ptrepack` that allows to (recursively) copy sub-hierarchies of objects present in a file into another one, changing, if desired, some of the filters applied to the leaves during the copy process.

Normally, these utilities will be installed somewhere in your `PATH` during the process of installation of the PyTables package, so that you can invoke them from any place in your file system after the installation has successfully finished.

3.4.1 ptdump

As has been said before, `ptdump` utility allows you look into the contents of your PyTables files. It lets you see not only the data but also the metadata (that is, the *structure* and additional information in the form of *attributes*).

Usage

For instructions on how to use it, just pass the `-h` flag to the command:

```
$ ptdump -h
```

to see the message usage:

```
usage: ptdump [-h] [-v] [-d] [-a] [-s] [-c] [-i] [-R RANGE]
            filename[:nodepath]
```

The ptdump utility allows you look into the contents of your PyTables files. It lets you see not only the data but also the metadata (that is, the **structure** and additional information in the form of **attributes**).

positional arguments:

```
filename[:nodepath]  name of the HDF5 file to dump
```

optional arguments:

```
-h, --help            show this help message and exit
-v, --verbose         dump more metainformation on nodes
-d, --dump            dump data information on leaves
-a, --showattrs       show attributes in nodes (only useful when -v or -d
                      are active)
-s, --sort            sort output by node name
-c, --colinfo         show info of columns in tables (only useful when -v or
                      -d are active)
-i, --idxinfo         show info of indexed columns (only useful when -v or
                      -d are active)
-R RANGE, --range RANGE
                      select a RANGE of rows (in the form "start,stop,step")
                      during the copy of *all* the leaves. Default values
                      are "None,None,1", which means a copy of all the rows.
```

Read on for a brief introduction to this utility.

A small tutorial on ptdump

Let's suppose that we want to know only the *structure* of a file. In order to do that, just don't pass any flag, just the file as parameter.

```
$ ptdump vllarray1.h5
/ (RootGroup) ''
/vllarray1 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of ints'
/vllarray2 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of strings'
```

we can see that the file contains just a leaf object called vllarray1, that is an instance of VLLArray, has 4 rows, and two filters has been used in order to create it: shuffle and zlib (with a compression level of 1).

Let's say we want more meta-information. Just add the -v (verbose) flag:

```
$ ptdump -v vllarray1.h5
/ (RootGroup) ''
/vllarray1 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of ints'
  atom = Int32Atom(shape=(), dflt=0)
  byteorder = 'little'
  nrow = 3
  flavor = 'numpy'
/vllarray2 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of strings'
  atom = StringAtom(itemsize=2, shape=(), dflt='')
```

(continues on next page)

(continued from previous page)

```
byteorder = 'irrelevant'
nrows = 3
flavor = 'python'
```

so we can see more info about the atoms that are the components of the `vlarray1` dataset, i.e. they are scalars of type `Int32` and with NumPy *flavor*.

If we want information about the attributes on the nodes, we must add the `-a` flag:

```
$ ptdump -va vlarray1.h5
/ (RootGroup) ''
/_v_attrs (AttributeSet), 4 attributes:
[CLASS := 'GROUP',
 PYTABLES_FORMAT_VERSION := '2.0',
 TITLE := '',
 VERSION := '1.0']
/vlarray1 (VLAarray(3,), shuffle, zlib(1)) 'ragged array of ints'
atom = Int32Atom(shape=(), dflt=0)
byteorder = 'little'
nrows = 3
flavor = 'numpy'
/vlarray1._v_attrs (AttributeSet), 3 attributes:
[CLASS := 'VLARRAY',
 TITLE := 'ragged array of ints',
 VERSION := '1.3']
/vlarray2 (VLAarray(3,), shuffle, zlib(1)) 'ragged array of strings'
atom = StringAtom(itemsizes=2, shape=(), dflt='')
byteorder = 'irrelevant'
nrows = 3
flavor = 'python'
/vlarray2._v_attrs (AttributeSet), 4 attributes:
[CLASS := 'VLARRAY',
 FLAVOR := 'python',
 TITLE := 'ragged array of strings',
 VERSION := '1.3']
```

Let's have a look at the real data:

```
$ ptdump -d vlarray1.h5
/ (RootGroup) ''
/vlarray1 (VLAarray(3,), shuffle, zlib(1)) 'ragged array of ints'
Data dump:
[0] [5 6]
[1] [5 6 7]
[2] [5 6 9 8]
/vlarray2 (VLAarray(3,), shuffle, zlib(1)) 'ragged array of strings'
Data dump:
[0] ['5', '66']
[1] ['5', '6', '77']
[2] ['5', '6', '9', '88']
```

We see here a data dump of the 4 rows in `vlarray1` object, in the form of a list. Because the object is a VLA, we see a different number of integers on each row.

Say that we are interested only on a specific *row range* of the `vlarray1` object:

```
ptdump -R2,3 -d vllarray1.h5:/vllarray1
/vllarray1 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of ints'
  Data dump:
[2] [5 6 9 8]
```

Here, we have specified the range of rows between 2 and 4 (the upper limit excluded, as usual in Python). See how we have selected only the /vllarray1 object for doing the dump (vllarray1.h5:/vllarray1).

Finally, you can mix several information at once:

```
$ ptdump -R2,3 -vad vllarray1.h5:/vllarray1
/vllarray1 (VLLArray(3,), shuffle, zlib(1)) 'ragged array of ints'
  atom = Int32Atom(shape=(), dtype=0)
  byteorder = 'little'
  nrows = 3
  flavor = 'numpy'
/vllarray1._v_attrs (AttributeSet), 3 attributes:
  [CLASS := 'VLLARRAY',
   TITLE := 'ragged array of ints',
   VERSION := '1.3']
  Data dump:
[2] [5 6 9 8]
```

3.4.2 ptrepack

This utility is a very powerful one and lets you copy any leaf, group or complete subtree into another file. During the copy process you are allowed to change the filter properties if you want so. Also, in the case of duplicated pathnames, you can decide if you want to overwrite already existing nodes on the destination file. Generally speaking, ptrepack can be useful in many situations, like replicating a subtree in another file, change the filters in objects and see how affect this to the compression degree or I/O performance, consolidating specific data in repositories or even *importing* generic HDF5 files and create true PyTables counterparts.

Usage

For instructions on how to use it, just pass the -h flag to the command:

```
$ ptrepack -h
```

to see the message usage:

```
usage: ptrepack [-h] [-v] [-o] [-R RANGE] [--non-recursive]
               [--dest-title TITLE] [--dont-create-sysattrs]
               [--dont-copy-userattrs] [--overwrite-nodes]
               [--complevel COMLEVEL]
               [--complib {zlib,lzo,bzip2,blosc,blosc:blosclz,blosc:lz4,blosc:lz4hc,
→blosc:snappy,blosc:zlib,blosc:zstd}]
               [--shuffle {0,1}] [--bitshuffle {0,1}] [--fletcher32 {0,1}]
               [--keep-source-filters] [--chunkshape CHUNKSHAPE]
               [--upgrade-flavors] [--dont-regenerate-old-indexes]
               [--sortby COLUMN] [--checkCSI] [--propindexes]
               sourcefile:sourcegroup destfile:destgroup
```

This utility is very powerful and lets you copy any leaf, group or complete subtree into another file. During the copy process you are allowed to change

(continues on next page)

(continued from previous page)

the filter properties **if** you want so. Also, in the **case** of duplicated pathnames, you can decide **if** you want to overwrite already existing nodes on the destination file. Generally speaking, ptrepack can be useful in many situations, like replicating a subtree in another file, change the filters in objects and see how affect this to the compression degree or I/O performance, consolidating specific data in repositories or even **importing** generic HDF5 files and create **true** PyTables counterparts.

positional arguments:

```

sourcefile:sourcegroup      source file/group
destfile:destgroup          destination file/group

```

optional arguments:

```

-h, --help                  show this help message and exit
-v, --verbose                show verbose information
-o, --overwrite              overwrite destination file
-R RANGE, --range RANGE     select a RANGE of rows (in the form "start,stop,step")
                              during the copy of *all* the leaves. Default values
                              are "None,None,1", which means a copy of all the rows.
--non-recursive              do not do a recursive copy. Default is to do it
--dest-title TITLE           title for the new file (if not specified, the source
                              is copied)
--dont-create-sysattrs       do not create sys attrs (default is to do it)
--dont-copy-userattrs        do not copy the user attrs (default is to do it)
--overwrite-nodes            overwrite destination nodes if they exist. Default is
                              to not overwrite them
--complevel COMLEVEL         set a compression level (0 for no compression, which
                              is the default)
--complib {zlib,lzo,bzip2,blosc,blosc:blosclz,blosc:lz4,blosc:lz4hc,blosc:snappy,
->blosc:zlib,blosc:zstd}      set the compression library to be used during the
                              copy. Defaults to zlib
--shuffle {0,1}              activate or not the shuffle filter (default is active
                              if complevel > 0)
--bitshuffle {0,1}           activate or not the bitshuffle filter (not active by
                              default)
--fletcher32 {0,1}           whether to activate or not the fletcher32 filter (not
                              active by default)
--keep-source-filters         use the original filters in source files. The default
                              is not doing that if any of --complevel, --complib,
                              --shuffle --bitshuffle or --fletcher32 option is
                              specified
--chunkshape CHUNKSHAPE      set a chunkshape. Possible options are: "keep" |
                              "auto" | int | tuple. A value of "auto" computes a
                              sensible value for the chunkshape of the leaves
                              copied. The default is to "keep" the original value
--upgrade-flavors             when repacking PyTables 1.x or PyTables 2.x files, the
                              flavor of leaves will be unset. With this, such a
                              leaves will be serialized as objects with the internal
                              flavor ('numpy' for 3.x series)

```

(continues on next page)

(continued from previous page)

```

--dont-regenerate-old-indexes    disable regenerating old indexes. The default is to
                                regenerate old indexes as they are found
--sortby COLUMN                 do a table copy sorted by the index in "column". For
                                reversing the order, use a negative value in the
                                "step" part of "RANGE" (see "-r" flag). Only applies
                                to table objects
--checkCSI                      force the check for a CSI index for the --sortby
                                column
--propindexes                   propagate the indexes existing in original tables. The
                                default is to not propagate them. Only applies to
                                table objects
--dont-allow-padding            remove the possible padding in compound types in
                                source files. The default is to propagate it. Only
                                applies to table objects

```

Read on for a brief introduction to this utility.

A small tutorial on ptrepack

Imagine that we have ended the tutorial 1 (see the output of `examples/tutorial1-1.py`), and we want to copy our reduced data (i.e. those datasets that hangs from the `/column` group) to another file. First, let's remember the content of the `examples/tutorial1.h5`:

```

$ ptdump tutorial1.h5
/ (RootGroup) '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'

```

Now, copy the `/columns` to other non-existing file. That's easy:

```
$ ptrepack tutorial1.h5:/columns reduced.h5
```

That's all. Let's see the contents of the newly created `reduced.h5` file:

```

$ ptdump reduced.h5
/ (RootGroup) ''
/name (Array(3,)) 'Name column selection'
/pressure (Array(3,)) 'Pressure column selection'

```

so, you have copied the children of `/columns` group into the *root* of the `reduced.h5` file.

Now, you suddenly realized that what you intended to do was to copy all the hierarchy, the group `/columns` itself included. You can do that by just specifying the destination group:

```

$ ptrepack tutorial1.h5:/columns reduced.h5:/columns
$ ptdump reduced.h5
/ (RootGroup) ''
/name (Array(3,)) 'Name column selection'
/pressure (Array(3,)) 'Pressure column selection'
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'

```

OK. Much better. But you want to get rid of the existing nodes on the new file. You can achieve this by adding the `-o` flag:

```
$ ptrepack -o tutorial1.h5:/columns reduced.h5:/columns
$ ptdump reduced.h5
/ (RootGroup) ''
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
```

where you can see how the old contents of the `reduced.h5` file has been overwritten.

You can copy just one single node in the repacking operation and change its name in destination:

```
$ ptrepack tutorial1.h5:/detector/readout reduced.h5:/rawdata
$ ptdump reduced.h5
/ (RootGroup) ''
/rawdata (Table(10,)) 'Readout example'
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
```

where the `/detector/readout` has been copied to `/rawdata` in destination.

We can change the filter properties as well:

```
$ ptrepack --complevel=1 tutorial1.h5:/detector/readout reduced.h5:/rawdata
Problems doing the copy from 'tutorial1.h5:/detector/readout' to 'reduced.h5:/rawdata'
The error was --> tables.exceptions.NodeError: destination group ``/`` already has
↳ a node named ``rawdata``; you may want to use the ``overwrite`` argument
The destination file looks like:
/ (RootGroup) ''
/rawdata (Table(10,)) 'Readout example'
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
Traceback (most recent call last):
  File "utils/ptrepack", line 3, in ?
    main()
  File ".../tables/scripts/ptrepack.py", line 349, in main
    stats = stats, start = start, stop = stop, step = step)
  File ".../tables/scripts/ptrepack.py", line 107, in copy_leaf
    raise RuntimeError, "Please check that the node names are not
    duplicated in destination, and if so, add the --overwrite-nodes flag
    if desired."
RuntimeError: Please check that the node names are not duplicated in
destination, and if so, add the --overwrite-nodes flag if desired.
```

Oops! We ran into problems: we forgot that the `/rawdata` pathname already existed in destination file. Let's add the `--overwrite-nodes`, as the verbose error message suggested:

```
$ ptrepack --overwrite-nodes --complevel=1 tutorial1.h5:/detector/readout
reduced.h5:/rawdata
$ ptdump reduced.h5
/ (RootGroup) ''
/rawdata (Table(10,), shuffle, zlib(1)) 'Readout example'
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
```

you can check how the filter properties has been changed for the /rawdata table. Check as the other nodes still exists.

Finally, let's copy a *slice* of the readout table in origin to destination, under a new group called /slices and with the name, for example, aslice:

```
$ ptrepack -R1,8,3 tutorial1.h5:/detector/readout reduced.h5:/slices/aslice
$ ptdump reduced.h5
/ (RootGroup) ''
/rawdata (Table(10,)), shuffle, zlib(1)) 'Readout example'
/columns (Group) ''
/columns/name (Array(3,)) 'Name column selection'
/columns/pressure (Array(3,)) 'Pressure column selection'
/slices (Group) ''
/slices/aslice (Table(3,)) 'Readout example'
```

note how only 3 rows of the original readout table has been copied to the new aslice destination. Note as well how the previously nonexistent slices group has been created in the same operation.

3.4.3 pt2to3

The PyTables 3.x series now follows [PEP 8](#) coding standard. This makes using PyTables more idiomatic with surrounding Python code that also adheres to this standard. The primary way that the 2.x series was *not* PEP 8 compliant was with respect to variable naming conventions. Approximately 450 API variables were identified and updated for PyTables 3.x.

To ease migration, PyTables ships with a new `pt2to3` command line tool. This tool will run over a file and replace any instances of the old variable names with the 3.x version of the name. This tool covers the overwhelming majority of cases was used to transition the PyTables code base itself! However, it may also accidentally also pick up variable names in 3rd party codes that have *exactly* the same name as a PyTables' variable. This is because `pt2to3` was implemented using regular expressions rather than a fancier AST-based method. By using regexes, `pt2to3` works on Python and Cython code.

`pt2to3` **help:**

```
usage: pt2to3 [-h] [-r] [-p] [-o OUTPUT] [-i] filename

PyTables 2.x -> 3.x API transition tool This tool displays to standard out, so
it is common to pipe this to another file: $ pt2to3 oldfile.py > newfile.py

positional arguments:
  filename              path to input file.

optional arguments:
  -h, --help            show this help message and exit
  -r, --reverse          reverts changes, going from 3.x -> 2.x.
  -p, --no-ignore-previous
                        ignores previous_api() calls.
  -o OUTPUT             output file to write to.
  -i, --inplace          overwrites the file in-place.
```

Note that `pt2to3` only works on a single file, not a a directory. However, a simple BASH script may be written to run `pt2to3` over an entire directory and all sub-directories:

```
#!/bin/bash
for f in $(find .)
do
    echo $f
```

(continues on next page)

(continued from previous page)

```
pt2to3 $f > temp.txt
mv temp.txt $f
done
```

Note: `pt2to3` uses the `argparse` module that is part of the Python standard library since Python 2.7. Users of Python 2.6 should install `argparse` separately (e.g. via `pip`).

3.5 PyTables File Format

PyTables has a powerful capability to deal with native HDF5 files created with another tools. However, there are situations where you may want to create truly native PyTables files with those tools while retaining full compatibility with PyTables format. That is perfectly possible, and in this appendix is presented the format that you should endow to your own-generated files in order to get a fully PyTables compatible file.

We are going to describe the *2.0 version of PyTables file format* (introduced in PyTables version 2.0). As time goes by, some changes might be introduced (and documented here) in order to cope with new necessities. However, the changes will be carefully pondered so as to ensure backward compatibility whenever is possible.

A PyTables file is composed with arbitrarily large amounts of HDF5 groups (Groups in PyTables naming scheme) and datasets (Leaves in PyTables naming scheme). For groups, the only requirements are that they must have some *system attributes* available. By convention, system attributes in PyTables are written in upper case, and user attributes in lower case but this is not enforced by the software. In the case of datasets, besides the mandatory system attributes, some conditions are further needed in their storage layout, as well as in the datatypes used in there, as we will see shortly.

As a final remark, you can use any filter as you want to create a PyTables file, provided that the filter is a standard one in HDF5, like *zlib*, *shuffle* or *gzip* (although the last one can not be used from within PyTables to create a new file, datasets compressed with *gzip* can be read, because it is the HDF5 library which does the decompression transparently).

3.5.1 Mandatory attributes for a File

The File object is, in fact, an special HDF5 *group* structure that is *root* for the rest of the objects on the object tree. The next attributes are mandatory for the HDF5 *root group* structure in PyTables files:

- **CLASS:** This attribute should always be set to 'GROUP' for group structures.
- **PYTABLES_FORMAT_VERSION:** It represents the internal format version, and currently should be set to the '2.0' string.
- **TITLE:** A string where the user can put some description on what is this group used for.
- **VERSION:** Should contain the string '1.0'.

3.5.2 Mandatory attributes for a Group

The next attributes are mandatory for *group* structures:

- **CLASS:** This attribute should always be set to 'GROUP' for group structures.
- **TITLE:** A string where the user can put some description on what is this group used for.
- **VERSION:** Should contain the string '1.0'.

3.5.3 Optional attributes for a Group

The next attributes are optional for *group* structures:

- **FILTERS**: When present, this attribute contains the filter properties (a *Filters* instance, see section [The Filters class](#)) that may be inherited by leaves or groups created immediately under this group. This is a packed 64-bit integer structure, where
 - *byte 0* (the least-significant byte) is the compression level (*complevel*).
 - *byte 1* is the compression library used (*complib*): 0 when irrelevant, 1 for Zlib, 2 for LZO and 3 for Bzip2.
 - *byte 2* indicates which parameterless filters are enabled (*shuffle* and *fletcher32*): bit 0 is for *Shuffle* while bit 1 is for *Fletcher32*.
 - other bytes are reserved for future use.

3.5.4 Mandatory attributes, storage layout and supported data types for Leaves

This depends on the kind of Leaf. The format for each type follows.

Table format

Mandatory attributes

The next attributes are mandatory for *table* structures:

- **CLASS**: Must be set to 'TABLE'.
- **TITLE**: A string where the user can put some description on what is this dataset used for.
- **VERSION**: Should contain the string '2.6'.
- **FIELD_X_NAME**: It contains the names of the different fields. The X means the number of the field, zero-based (beware, order do matter). You should add as many attributes of this kind as fields you have in your records.
- **FIELD_X_FILL**: It contains the default values of the different fields. All the datatypes are supported natively, except for complex types that are currently serialized using Pickle. The X means the number of the field, zero-based (beware, order do matter). You should add as many attributes of this kind as fields you have in your records. These fields are meant for saving the default values persistently and their existence is optional.
- **NROWS**: This should contain the number of *compound* data type entries in the dataset. It must be an *int* data type.

Storage Layout

A Table has a *dataspace* with a *1-dimensional chunked* layout.

Datatypes supported

The datatype of the elements (rows) of Table must be the `H5T_COMPOUND` *compound* data type, and each of these compound components must be built with only the next HDF5 data types *classes*:

- **H5T_BITFIELD**: This class is used to represent the Bool type. Such a type must be build using a `H5T_NATIVE_B8` datatype, followed by a `HDF5 H5Tset_precision` call to set its precision to be just 1 bit.
- **H5T_INTEGER**: This includes the next data types:

- *H5T_NATIVE_SCHAR*: This represents a *signed char* C type, but it is effectively used to represent an Int8 type.
 - *H5T_NATIVE_UCHAR*: This represents an *unsigned char* C type, but it is effectively used to represent an UInt8 type.
 - *H5T_NATIVE_SHORT*: This represents a *short* C type, and it is effectively used to represent an Int16 type.
 - *H5T_NATIVE_USHORT*: This represents an *unsigned short* C type, and it is effectively used to represent an UInt16 type.
 - *H5T_NATIVE_INT*: This represents an *int* C type, and it is effectively used to represent an Int32 type.
 - *H5T_NATIVE_UINT*: This represents an *unsigned int* C type, and it is effectively used to represent an UInt32 type.
 - *H5T_NATIVE_LONG*: This represents a *long* C type, and it is effectively used to represent an Int32 or an Int64, depending on whether you are running a 32-bit or 64-bit architecture.
 - *H5T_NATIVE_ULONG*: This represents an *unsigned long* C type, and it is effectively used to represent an UInt32 or an UInt64, depending on whether you are running a 32-bit or 64-bit architecture.
 - *H5T_NATIVE_LLONG*: This represents a *long long* C type (`__int64`, if you are using a Windows system) and it is effectively used to represent an Int64 type.
 - *H5T_NATIVE_ULLONG*: This represents an *unsigned long long* C type (beware: this type does not have a correspondence on Windows systems) and it is effectively used to represent an UInt64 type.
- ***H5T_FLOAT*: This includes the next datatypes:**
 - *H5T_NATIVE_FLOAT*: This represents a *float* C type and it is effectively used to represent an Float32 type.
 - *H5T_NATIVE_DOUBLE*: This represents a *double* C type and it is effectively used to represent an Float64 type.
 - ***H5T_TIME*: This includes the next datatypes:**
 - *H5T_UNIX_D32*: This represents a POSIX *time_t* C type and it is effectively used to represent a ‘Time32’ aliasing type, which corresponds to an Int32 type.
 - *H5T_UNIX_D64*: This represents a POSIX *struct timeval* C type and it is effectively used to represent a ‘Time64’ aliasing type, which corresponds to a Float64 type.
 - *H5T_STRING*: The datatype used to describe strings in PyTables is *H5T_C_S1* (i.e. a *string* C type) followed with a call to the HDF5 *H5Tset_size()* function to set their length.
 - *H5T_ARRAY*: This allows the construction of homogeneous, multidimensional arrays, so that you can include such objects in compound records. The types supported as elements of *H5T_ARRAY* data types are the ones described above. Currently, PyTables does not support nested *H5T_ARRAY* types.
 - *H5T_COMPOUND*: This allows the support for datatypes that are compounds of compounds (this is also known as *nested types* along this manual).

This support can also be used for defining complex numbers. Its format is described below:

The *H5T_COMPOUND* type class contains two members. Both members must have the *H5T_FLOAT* atomic datatype class. The name of the first member should be “r” and represents the real part. The name of the second member should be “i” and represents the imaginary part. The *precision* property of both of the *H5T_FLOAT* members must be either 32 significant bits (e.g. *H5T_NATIVE_FLOAT*) or 64 significant bits (e.g. *H5T_NATIVE_DOUBLE*). They represent *Complex32* and *Complex64* types respectively.

Array format

Mandatory attributes

The next attributes are mandatory for *array* structures:

- *CLASS*: Must be set to 'ARRAY'.
- *TITLE*: A string where the user can put some description on what is this dataset used for.
- *VERSION*: Should contain the string '2.3'.

Storage Layout

An Array has a *dataspace* with a *N-dimensional contiguous* layout (if you prefer a *chunked* layout see EArray below).

Datatypes supported

The elements of Array must have either HDF5 *atomic* data types or a *compound* data type representing a complex number. The atomic data types can currently be one of the next HDF5 data type *classes*: H5T_BITFIELD, H5T_INTEGER, H5T_FLOAT and H5T_STRING. The H5T_TIME class is also supported for reading existing Array objects, but not for creating them. See the Table format description in [Table format](#) for more info about these types.

In addition to the HDF5 atomic data types, the Array format supports complex numbers with the H5T_COMPOUND data type class. See the Table format description in [Table format](#) for more info about this special type.

You should note that H5T_ARRAY class datatypes are not allowed in Array objects.

CArray format

Mandatory attributes

The next attributes are mandatory for *CArray* structures:

- *CLASS*: Must be set to 'CARRAY'.
- *TITLE*: A string where the user can put some description on what is this dataset used for.
- *VERSION*: Should contain the string '1.0'.

Storage Layout

An CArray has a *dataspace* with a *N-dimensional chunked* layout.

Datatypes supported

The elements of CArray must have either HDF5 *atomic* data types or a *compound* data type representing a complex number. The atomic data types can currently be one of the next HDF5 data type *classes*: H5T_BITFIELD, H5T_INTEGER, H5T_FLOAT and H5T_STRING. The H5T_TIME class is also supported for reading existing CArray objects, but not for creating them. See the Table format description in [Table format](#) for more info about these types.

In addition to the HDF5 atomic data types, the CArray format supports complex numbers with the H5T_COMPOUND data type class. See the Table format description in [Table format](#) for more info about this special type.

You should note that H5T_ARRAY class datatypes are not allowed yet in Array objects.

EArray format

Mandatory attributes

The next attributes are mandatory for *earray* structures:

- **CLASS**: Must be set to 'EARRAY'.
- **EXTDIM**: (*Integer*) Must be set to the extendable dimension. Only one extendable dimension is supported right now.
- **TITLE**: A string where the user can put some description on what is this dataset used for.
- **VERSION**: Should contain the string '1.3'.

Storage Layout

An EArray has a *dataspace* with a *N-dimensional chunked* layout.

Datatypes supported

The elements of EArray are allowed to have the same data types as for the elements in the Array format. They can be one of the HDF5 *atomic* data type *classes*: H5T_BITFIELD, H5T_INTEGER, H5T_FLOAT, H5T_TIME or H5T_STRING, see the Table format description in [Table format](#) for more info about these types. They can also be a H5T_COMPOUND datatype representing a complex number, see the Table format description in [Table format](#).

You should note that H5T_ARRAY class data types are not allowed in EArray objects.

VArray format

Mandatory attributes

The next attributes are mandatory for *varray* structures:

- **CLASS**: Must be set to 'VLARRAY'.
- **PSEUDOATOM**: This is used so as to specify the kind of pseudo-atom (see [VArray format](#)) for the VArray. It can take the values 'vlstring', 'vlunicode' or 'object'. If your atom is not a pseudo-atom then you should not specify it.
- **TITLE**: A string where the user can put some description on what is this dataset used for.
- **VERSION**: Should contain the string '1.3'.

Storage Layout

An VArray has a *dataspace* with a *1-dimensional chunked* layout.

Data types supported

The data type of the elements (rows) of VArray objects must be the H5T_VLEN *variable-length* (or VL for short) datatype, and the base datatype specified for the VL datatype can be of any *atomic* HDF5 datatype that is listed in the Table format description *Table format*. That includes the classes:

- H5T_BITFIELD
- H5T_INTEGER
- H5T_FLOAT
- H5T_TIME
- H5T_STRING
- H5T_ARRAY

They can also be a H5T_COMPOUND data type representing a complex number, see the Table format description in *Table format* for a detailed description.

You should note that this does not include another VL datatype, or a compound datatype that does not fit the description of a complex number. Note as well that, for object and vlstring pseudo-atoms, the base for the VL datatype is always a H5T_NATIVE_UCHAR (H5T_NATIVE_UINT for vlunicode). That means that the complete row entry in the dataset has to be used in order to fully serialize the object or the variable length string.

3.5.5 Optional attributes for Leaves

The next attributes are optional for *leaves*:

- *FLAVOR*: This is meant to provide the information about the kind of object kept in the Leaf, i.e. when the dataset is read, it will be converted to the indicated flavor. It can take one the next string values:
 - “*numpy*”: Read data (structures arrays, arrays, records, scalars) will be returned as NumPy objects.
 - “*python*”: Read data will be returned as Python lists, tuples, or scalars.

3.6 Bibliography

[*HDFG1*] The HDF Group. What is HDF5?. Concise description about HDF5 capabilities and its differences from earlier versions (HDF4). <http://www.hdfgroup.org/HDF5/whatishdf5.html>.

[*HDFG2*] The HDF Group. Introduction to HDF5. Introduction to the HDF5 data model and programming model. <http://www.hdfgroup.org/HDF5/doc/H5.intro.html>.

[*HDFG3*] The HDF Group. The HDF5 table programming model. Examples on using HDF5 tables with the C API. <http://www.hdfgroup.org/HDF5/Tutor/h5table.html>.

[*MERTZ*] David Mertz. 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. http://gnosis.cx/publish/programming/xml_matters_2.html.

[*CYTHON*] Stefan Behnel, Robert Bradshaw, Dag Sverre Seljebotn, and Greg Ewing. Cython. A language that makes writing C extensions for the Python language as easy as Python itself. <http://www.cython.org>.

[*NUMPY*] Travis Oliphant and et al. NumPy. Scientific Computing with Numerical Python. The latest and most powerful re-implementation of Numeric to date. It implements all the features that can be found in Numeric and numarray, plus a bunch of new others. In general, it is more efficient as well. <http://www.numpy.org>.

- [NUMEXPR]** David Cooke, Francesc Altèd, and et al. Numexpr. Fast evaluation of array expressions by using a vector-based virtual machine. It is an enhanced computing kernel that is generally faster (between 1x and 10x, depending on the kind of operations) than NumPy at evaluating complex array expressions. <http://code.google.com/p/numexpr>.
- [ZLIB]** JeanLoup Gailly and Mark Adler. zlib. A Massively Spiffy Yet Delicately Unobtrusive Compression Library. A standard library for compression purposes. <http://www.gzip.org/zlib/>.
- [LZO]** Markus F Oberhumer. LZO. A data compression library which is suitable for data de-/compression in real-time. It offers pretty fast compression and decompression with reasonable compression ratio. <http://www.oberhumer.com/opensource/>.
- [BZIP2]** Julian Seward. bzip2. A high performance lossless compressor. It offers very high compression ratios within reasonable times. <http://www.bzip.org/>.
- [BLOSC]** Francesc Altèd. Blosc. A blocking, shuffling and loss-less compression library. A compressor designed to transmit data from memory to CPU (and back) faster than a plain memcpy(). <http://www.blosc.org/>.
- [GNUWIN32]** Alexis Wilke, Jerry S., Kees Zeelenberg, and Mathias Michaelis. GnuWin32. GNU (and other) tools ported to Win32. GnuWin32 provides native Win32-versions of GNU tools, or tools with a similar open source licence. <http://gnuwin32.sourceforge.net/>.
- [PSYCO]** Armin Rigo. Psyco. A Python specializing compiler. Run existing Python software faster, with no change in your source. <http://psyco.sourceforge.net>.
- [SCIPY1]** Konrad Hinsén. Scientific Python. Collection of Python modules useful for scientific computing. <http://dirac.cnrs-orleans.fr/ScientificPython>.
- [SCIPY2]** Eric Jones, Travis Oliphant, Pearu Peterson, and et al. SciPy. Scientific tools for Python. SciPy supplements the popular Numeric module, gathering a variety of high level science and engineering modules together as a single package. <http://www.scipy.org>.
- [OPTIM]** Francesc Altèd and Ivan Vilata. Optimization of file openings in PyTables. This document explores the savings of the opening process in terms of both CPU time and memory, due to the adoption of a LRU cache for the nodes in the object tree. <http://www.pytables.org/docs/NewObjectTreeCache.pdf>.
- [OPSI]** Francesc Altèd and Ivan Vilata. OPSI: The indexing system of PyTables 2 Professional Edition. Exhaustive description and benchmarks about the indexing engine that comes with PyTables Pro. <http://www.pytables.org/docs/OPSI-indexes.pdf>.
- [VITABLES]** Vicent Mas. ViTables. A GUI for PyTables/HDF5 files. It is a graphical tool for browsing and editing files in both PyTables and HDF5 formats. <http://vitables.org>.
- [GIT]** Git is a free and open source, distributed version control system designed to handle everything from small to very large projects with speed and efficiency <http://git-scm.com>.
- [SPHINX]** Sphinx is a tool that makes it easy to create intelligent and beautiful documentation, written by Georg Brandl and licensed under the BSD license <http://sphinx-doc.org>.

Symbols

`_v_pos` (*tables.Col attribute*), 64

A

`atom` (*tables.Array attribute*), 62

B

`BLOSC_DIR`, 12, 14, 15

`BZIP2_DIR`, 14, 15

E

environment variable

`BLOSC_DIR`, 12, 14, 15

`BZIP2_DIR`, 14, 15

`HDF5_DIR`, 14, 15

`LD_LIBRARY_PATH`, 14

`LIBS`, 14

`LZO_DIR`, 14, 15

`PATH`, 17

`PYTHONPATH`, 16–18

`USE-PKGCONFIG`, 14

H

`HDF5_DIR`, 14, 15

L

`LD_LIBRARY_PATH`, 14

`LIBS`, 14

`LZO_DIR`, 14, 15

N

`nelements` (*tables.tables.index.Index attribute*), 65

`nrow` (*tables.Array attribute*), 62

P

`PATH`, 17

`PYTHONPATH`, 16–18

S

`size_in_memory` (*tables.Leaf attribute*), 61

U

`USE-PKGCONFIG`, 14