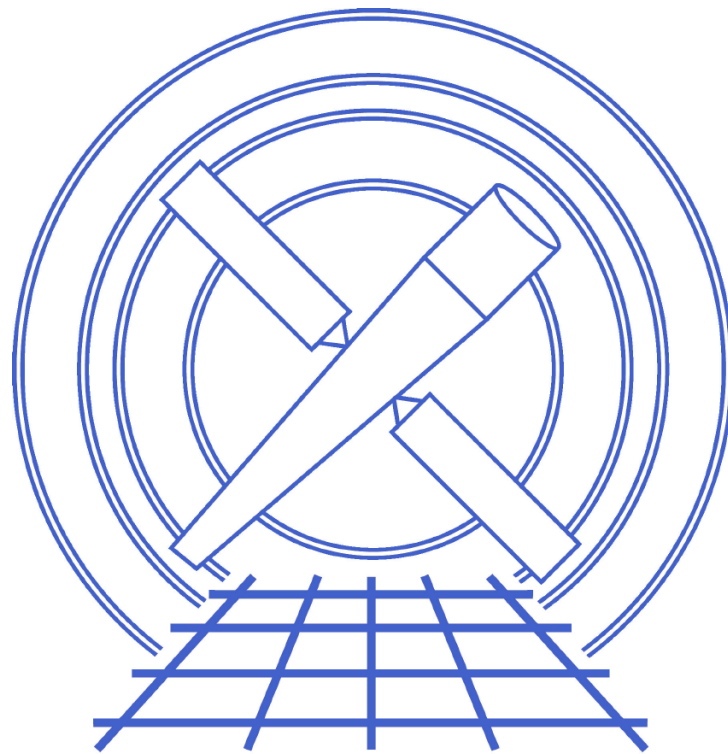


CXC-DM-002

CXC Data Model



Vol. 2

Data Model Abstract Design

Jonathan McDowell
Chandra X-ray Center
October 22, 2001

Contents

1	Introduction	4
1.1	Overview of the Model	4
1.2	What is a Data Model?	4
1.3	Summary of motivation	6
1.4	Problems and Solutions	6
1.5	Informal Introduction to the Data Model	9
1.6	Table columns	9
1.7	Table Attributes	10
1.8	Binned Data	10
1.9	Arrays and Images	10
1.10	Descriptors and Elements and Components	10
1.11	Stacks	11
2	Some general requirements	13
2.1	Data Model and files	13
2.2	Compatibility Requirements on FITS kernel	13
2.3	The native data model in FITS	13
2.4	Interaction of Data Model and other infrastructure	15
3	The CXC Data Model, SDS Version 2.0	15
3.1	DM Table	17
4	Table Data Section	17
4.1	Table Data	17
4.2	Descriptor	19
4.3	Array Dimensions	23
4.4	Array Axis	23
4.5	Axis Groups	23
4.6	Parent Descriptor	24
4.7	Coordinate Transform Descriptor	24
4.8	Column Data Descriptor	27
4.9	Interval type	27
4.10	Elements	28
4.11	Region Description	30
4.12	Table Data Cell	30
4.13	Table Row	30
4.14	Descriptor Groups	30
5	Data Subspace	31
5.1	Introduction	31
5.2	Unions of subspaces	32
5.3	General definition	33
6	Header	37
6.1	Key Data Descriptor	37
6.2	Grouping Descriptors	38

7 DM Images	38
7.1 Images and Tables	38
8 Case studies and examples	38
8.1 FITS case study: PSPC off axis histogram file	38
8.2 Case Study: Barycenter Correction Algorithm	49

1 Introduction

The Science Data Systems Group at the Chandra X-ray Center (CXC) developed the CXC Data Model (DM) as a generic data model for astronomical data, underlying the CIAO data analysis system and Chandra processing pipelines.

This document describes the design of version 2 of the DM developed in 2001.

1.1 Overview of the Model

Our data model has a number of high level goals:

- Create data files which are more fully self-describing, while retaining back compatibility in the sense that existing archival FITS files will be interpreted correctly.
- Systematize the treatment of data filtering, units, and coordinate systems, unifying the current approach which involves a large number of special cases.
- Allow programs to use both FITS and native IRAF file formats interchangeably, by supplying a format-independent interface layer.
- Allow users to write their own programs easily by providing a subroutine interface which makes accessing the data easy and removes the need for the user to worry about the details of the file format.
- Support advanced virtual file and filtering operations by providing a uniform convention for recording the way a file has been filtered.

The intent of the model is to describe an abstract representation of a generic astronomical dataset and to layer extra structure onto existing file formats to make them more fully self-defining. Our datasets include both binned (image) and tabular data, corresponding to the IRAF IMH and QPOE formats or the FITS IMAGE and BINTABLE formats. In this document I will make explicit parallels to the FITS format since it is the externally defined export and exchange format. An important aspect of the design presented here is that existing FITS files will be interpreted correctly by the data model. This is achieved through careful use of default values for keywords in the mapping to FITS.

1.2 What is a Data Model?

A data model is an abstract description of our datasets. (Datasets may be files, or groups of files that we want to consider as a unit). It tells us the different properties and attributes a dataset can have (e.g. 'a dataset consists of a header and a table or an array; a table has n columns each with a name, a data type, a unit, ' ... etc.) This description of the data is possible because all of our many different datasets can be thought of as special cases of a very small number of basic types of dataset. In using the data model to describe a dataset, then, we have a way of defining that dataset which makes explicit its differences from all other datasets. Furthermore, the data model contains no information about the storage format of the data. Thus our definition of the structure of a dataset is completely separated from the way in which that structure is implemented on disk - we distinguish between information that is truly part of the scientific data and information that is bookkeeping or specific to the file format. This makes it easy to support multiple file formats with the same data model. The data model can be implemented as an API which lets you access and manipulate the data using the concepts of the data model.

The data model gives the application writer an interface to the data which is independent of the details of the file format. It also provides a standardized structure and language which brings out the similarities between different kinds of dataset. This standardization is an important advance beyond the standardization provided by particular data formats such as FITS.

- We make the treatment of the data independent of the choice of disk file data format, thus allowing the algorithm to concentrate on the science and making it easy to support the open architecture of different data formats. It means that applications writers don't have to worry about the specifics of the data format, those are hidden in the interface subroutines.
- The model layers extra structure onto the concepts implicit in the underlying data formats.
- We describe all data in a common structure; by imposing a uniform description we can support generic tools. We have a way of describing a general data file independent of the specific structure of the file (PHA file, event file, etc). This means that when you make a new kind of file, existing tools can still do something with it. FTOOLS does this at a certain level, allowing basic filtering of generic tables, and can be thought of as having a very simple data model consisting of table columns with no extra attributes. Our software will go well beyond this, dealing with coordinate systems and other auxiliary quantities in a standardized way.
- Further, all this makes the data more self-describing.
- We explicitly tie information relating to each image axis or table column to that axis or column. In FITS, there is some of this: a keyword like CRVAL4 tells you the coordinate value for axis number 4. However, there are a lot of other keywords that don't do this and could - for instance, TSTART gives the start time for a dataset, but there is no explicit expression of the fact that this quantity is related to the TIME column in the data. The FITS kernel to the DM understands this particular association, and puts it into the DM structure which allows a generic association of a range of values with a column.
- Note that a single data model table may correspond to many FITS tables. For instance, the Good Time Intervals, which in the data model are just the ranges for one axis of the data subspace, have to be kept in a separate table in the FITS file. At the moment FITS files often have an assortment of tables in them, some of which are related to each other and some of which aren't. Using a data model helps us make much more sensible decisions about which FITS tables to group together in a single file. For instance, for an EVENT file it helps us realize that the Good Time Intervals are truly just an auxiliary piece of information describing the main table, while ROSAT Temporal Status Intervals are (at least on the data model I present here) a separate data object that has meaning separately from the EVENT data.
- The concept of a data subspace lets us unify the treatment of good time intervals, spatial regions, and filter ranges. This makes these concepts independent of whether a particular column contains temporal, spatial or spectral info, and lets us be much more systematic about asking the question 'to what range of data values does this dataset apply?'
- Grouping together of header keywords helps us propagate related info more easily, makes it easier to specify the definition of new files in terms of old ones, and improves user readability of headers.
- The existence of a data model helps us include support for new features (e.g. uncertainties) in a systematic way, so we don't have to deal with hundreds of special cases each time. This applies both to the new features we add now and to future features in later versions of the model - in other words, having an overall data model reduces overhead in including new functionality, because it's clear how to add that new functionality in a way that will work throughout the system.
- The separates out the science description from the details of a data format, allowing us to define clean mappings to different data formats. This makes it easier to support new data formats, since the I/O is so well isolated.

How can we be all things to all systems? The crucial idea is the concept of a data model. By this we don't mean a model of a specific dataset, like a spectral fitting functional model, we mean a model of the *concept* of

astronomical data. More specifically, we mean an abstract description of the structure of our data separate from its implementation in a particular disk storage format. We note for the software-literate that this abstract description can be - but does not need to be - given a manifest software implementation as an object or set of objects in an object-oriented language. Once we have our data model, we can map it to the particular disk data formats we wish to support. This allows the same code to read FTOOLS FITS files or PROS QPOE files and ‘see’ them (after a translation layer) as identical sources of information. The individual tool will not usually need any explicit ‘if FITS then’ code, and will not even know what type of file is being read.

In principle such a data model could be arbitrarily complex with many special cases. Actually it turns out that almost all our data can be described by a single kind of object, perhaps with a few simple flavors. This fact is what gives FTOOLS its strength: much X-ray data analysis can be accomplished by fairly general manipulations of FITS binary tables. We take FTOOLS’ advance one step farther by separating our unified data description from the specifics of the FITS format (2880-byte blocks, indexed keywords, storing the structure of the main data as header keywords, no units on keywords, etc), which are not relevant to any of the science algorithms. This separation turns out to be extremely powerful, and allows us to do a lot more than just support multiple data analysis contexts.

The existing package tool kits (FDUMP, TPRINT, etc..) will work on our files but may lose the extra layers of meaning provided by our data model. We therefore provide a new set of infrastructure tools which will do generic operations on our files. We have unified and extended the PROS concepts of filters, regions and good time intervals into a single selector concept; this greatly increases the flexibility of filtering.

1.3 Summary of motivation

By generalizing our approach, we can get by with fewer distinct tools. By writing the tools using our data model, and modern software approaches (careful layering, self-describing data, etc.), we can make each tool more flexible, able to do sensible things with data that is in slightly different formats, or even data representing entirely different physical quantities. We try and strip the algorithm to its bare bones and encode the specifics of the data in the self describing data files, not in the compiled code. By designing in low level support for operations on multiple data files, we make easier the task of doing the same operation across such sets of data files and, if desired, combining the results. The existence of a unified data model makes communications among programs, and between programs and GUIs, easier to systematize. Eventually, by including uncertainties, upper limits, units, etc., in our data model, we will standardize their treatment and so allow generic tools to operate on them.

1.4 Problems and Solutions

In this section I discuss various limitations we’ve come across in the way current systems handle abstract data manipulation. I concentrate on examples from PROS and FITS since they are the systems I am most familiar with.

- PROBLEM: PROS regions are handled in a different way from time, PHA filters.
- SOLUTION: Introduce the idea of a Data Subspace which handles filters on all data axes in a uniform way. The user can specify a spatial region anywhere they can specify a PHA or time filter. The Data Subspace for a data object records the way that object has been filtered. If you like, it is the filter that has been applied to the data so far. The Good Time Intervals are part of this filter.
- PROBLEM: Making a detector coordinate image was messy in PROS (PROS keyx, keyy syntax).

```
qplist "test.qp[pi=40:90]" region="c 2048 2048 20"
```

lists photons in a given sky region and PI range, but

```
display "test.qp[pi=40:90]"
```

does not take a region argument - you can't display it. To list photons in a detector coordinate region,

```
qplist "test.qp[key=(detx,dety),pi=40:90]"  
      region="c 2048 2048 20"
```

which is ugly because the specification of the region and the statement that the region applies to detector coordinates are separated.

- SOLUTION(1): Make regions part of the virtual file syntax, so you can do:

```
dmlist "test.qp[(detx,dety)=circle(2048,2048,20),pi=40:90]"
```

- this is much more coherent.

- SOLUTION(2) The scientist thinks in terms of 'detector position' and 'sky position' as single attributes of the data. Make our software able to work on two-dimensional items named 'DET' and 'SKY' to allow a natural system of

```
dmlist "test.qp[det=circle(2048,2048,20),pi=40:90]"
```

Make the data model support 2D objects with a name for the object and for each of its components (e.g. object name SKY, component names RA and DEC). This makes it easy for a programmer to make a file which knows that it contains a bunch of SKY each of which consists of an RA and a DEC. Current files don't have any way of letting the software know which columns are paired together as positions.

- PROBLEM: No standard way to record how the data has been filtered on PHA or PI.
- SOLUTION: The Data Subspace does this automatically. Thus the software can, if properly coded, know where to look to find out which (energy-dependent) point spread function would be matched to the current image - it looks for a PI axis in the image's data subspace.
- PROBLEM: Some data manipulation tasks need you to go back and forth between header keywords and table columns, but header keywords in FITS don't contain as much information as table columns (short names, no units, no vectors). Examples: we wish to combine event lists from ACIS chips I2 and I3, which have header CHIP_ID values giving the chip ID, getting an event list with an extra CHIP_ID column in which each row is either I2 or I3. Or, we wish to combine tables of sources detected with three different cell sizes, to make one table with a CELL_SIZE column. The resulting table needs to know the units in which CELL_SIZE is measured. Actually, it would currently have to be CELL_SIZ since the header keywords can only be 8 characters.
- SOLUTION: The data model supports the extra information. The I/O library handles a convention to write this to FITS in a way that is back compatible with existing data, and has now been incorporated as well in Goddard's FITSIO software. The tool program can ask for the same information about a keyword that it would for a column entry, so the code is more uniform - fewer special cases.
- PROBLEM: We have a blocked sky image and want to know about both the original plane pixel coordinate system and the celestial spherical coordinate system. Most FITS applications only support one set of WCS keywords for an image.

- SOLUTION: For array objects, allow an ‘physical’ coordinate system and an ‘world’ coordinate system to retain both sets of information. Allow arbitrary numbers of world coordinate systems for each object, so that for instance one could attach a galactic coordinate system to the image as well. We use FITS keyword conventions that are under consideration for adding to the FITS standard.
- PROBLEM: Want a single program to browse and plot all kinds of data files, labelling axes sensibly.
- SOLUTION: Each axis in the data is liable to have both a local and a ‘world’ value: pixel position and celestial position, mission time in seconds and calendar date, pulse height and nominal energy. The data model treats all of these as generic coordinate systems, so a plotting program can recognize them automatically. Example: pulse height versus time image, with nominal energies and calendar dates automatically labelled.
- PROBLEM: Want to support a table with images embedded in one of the columns, for instance aspect camera records.
- SOLUTION: Introduce a data model convention to handle this case, which is supported to a limited extent in FITS by the multidimensional array TDIMn syntax; further simplify by considering an ordinary image to be a special case of a table with one row and column.
- PROBLEM: Want to create datasets such as an array of x-ray colors versus best fit parameters, and invert to make an array of best fit parameters versus x-ray colors.
- SOLUTION: Provide data model support for arrays whose elements are themselves n-dimensional.
- PROBLEM: In a derived file like a light curve, we may make many columns (raw counts, background counts, net count rate, etc.) even though the basic concept is of time versus net count rate. We want our plotting software to plot the two columns of most interest by default. Also, indexing operations may also be carried out on event list columns of ‘most interest’.
- SOLUTION: Define ‘preferred’ columns (axes) of the table, which will rank a subset of the columns in an order which may be different from the order of the columns in the table. A plotting program which plots two descriptors (quantities) against each other will then take the first two preferred columns if such exist, otherwise it will take the first two columns in the table column order.
- PROBLEM: Want to deal with upper limits properly.
- SOLUTION: The interface to the data files should be able to cope with any data item being either a detection, an upper limit, or a detection with uncertainty. Other software, however, will see the uncertainty ranges as separate columns and won’t know that a particular value is an upper limit.
This solution is not yet implemented in the CXC DM, but support for NaN entries and null entries is present.
- PROBLEM: We have a set of PSFs which were created at XRCF at different energies; they are labelled with a ‘header keyword’ ENERGY. We wish to plot the FWHM of the PSFs versus energy. In an existing system, one would run the calculate-FWHM program on each PSF file separately, capturing the results and running a table creation program to combine them in a single result table (or noting them down on paper and typing them back in!); plotting the results might not be trivial either.
- SOLUTION: We should be able to do this with three commands: one to stack the PSF files on energy, creating an index file consisting of a table of energy versus filename, a second to run the calculate-FWHM program on the stack, and the third to plot the resulting file. In our system, the added bonus is that if the calculate-FWHM program also calculates uncertainties, these will be picked up by the plotting program.
This solution is not yet supported; a special stacking tool (to make stacks with index columns like energy...) is needed to do this and will be considered for future development.

1.5 Informal Introduction to the Data Model

In our model, each dataset consists of an ordered set of ‘Blocks’. A DM Block consists of a table whose columns may be scalars, vectors, arrays, or ranges. Header descriptors may be attached to the table as a whole, or to individual columns or to the data subspace. An important special case of a block table is called an Image, and we will often consider blocks to be of two types, Table and Image (even though strictly speaking an Image can be treated as a kind of Table).

A table consists of a header, together with a set of rows and columns. I will refer to the intersection of a row and a column as a ‘cell’. Some of our tabular data products will contain small embedded images. For instance, aspect camera data will include 6x6 pixel images of each fiducial light in every row of the table (the row represents the information from one aspect camera exposure), and FAINT mode event data has 3x3 images of the event island in the PHAS column. Also, we may eventually include small ‘postage-stamp’ images of sources in our source list data product. This suggests a theoretical simplification: if an image can be in a cell of a table, we may consider an image on its own to be just a table with one cell. So, instead of two different fundamental types of data, we have a single type - the ‘table-which-can-contain-images’. We can then specify that a table with one row and one column may, if desired, be stored on disk using an image format instead of a binary table format. We have thus moved the distinction between a table file and an image file to a different level: an image is a component of a table, rather than its peer. We still, of course, need to have interfaces to operate on images, so this simplification is minor in practical terms.

The header in a FITS file is a heterogeneous collection of information. Some of the keywords describe the file’s structure, while the remainder are metadata: data which apply to the file as a whole, but are true science data rather than descriptive of the file structure. We want to layer extra structure on the file so we can tell the difference between these types of header keyword. Some of the metadata has particular importance: it describes how the data in the table columns was selected. We treat this kind of information in a systematic way and isolate it conceptually as the table’s ‘data subspace’.

The figure below gives a schematic example of a complicated table.

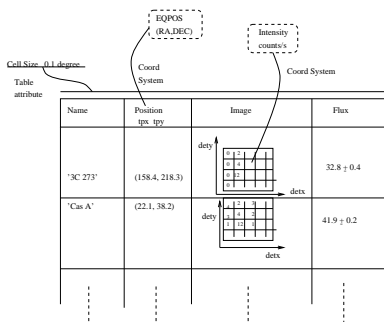


Figure 1: Example of a complicated table. The table is a source list containing ‘postage stamp’ images of each source. The position column has a coordinate system attached to it, the flux column has uncertainties, and the whole table has metadata such as the source detection cell size.

1.6 Table columns

FITS already provides support for vectors and arrays in table columns. However, there are several enhancements we need. Particularly for the case of positional data, we want to have paired table columns: for instance, DETX and DETY paired as DETPOS, or RA and DEC paired as EQPOS, with both the individual and the pair names available in the file. We also want to support uncertainties and upper limits, which implies something like having

a column FLUX and a column FLUX_ERR (no problem right now) together with a structure which ties the two together as a single object (Flux with error). Both of these enhancements, and the desire for back compatibility, lead us to a system with a low level (FITS) set of raw columns and a high level (Data Model) set of columns, with one high level column mapping to several low level columns.

1.7 Table Attributes

Table attributes are the equivalent of header keywords. Unlike FITS header keywords, we support the various descriptor attributes such as units, etc. FITS allows 'indexed keywords' which are really 1-D arrays of keywords: we want to support this at a higher level, and add support for 'vector keywords', e.g. grouping together RA and DEC as a single high level table attribute EQPOS.

We'd also like to specify some attributes as belonging to specific table columns rather than to the table as a whole. These are called column attributes. Similarly, the data subspace may have its own attributes: livetime is an example.

1.8 Binned Data

An event list table consists of values which represent precise points in an n-dimensional space. In contrast, we often deal with binned data in which the values represent cells of finite volume in the space. The simplest example is a histogram with equal size bins, but we also have datasets with logarithmic bins or even arbitrary bins (e.g. those chosen to match the position of sharp features in a spectrum). A binned data column can use the same mechanism as the uncertainties for a normal column, since it just involves specifying a range.

1.9 Arrays and Images

When we have binned data with ordered, equal size, contiguous bins, the column of data may be defined implicitly by specifying the start value and step size. Suppose we have a table whose columns include three binned data columns and two point data columns, one of which happens to be a 3D position:

C1	C2	C3	C4	C5
[0.5:1.5)	[10.0:11.0)	[4.8:4.9)	1082.2	(0.0, 18.3, -812.3)
[1.5:2.5)	[10.0:11.0)	[4.8:4.9)	182.3	(4.3, 12.2, -712.3)
....				
[0.5:1.5)	[11.0:12.0)	[4.9:5.0)	1211.2	(2.1, -1.2, -271.3)
[1.5:2.5)	[11.0:12.0)	[4.9:5.0)	1232.1	(6.2, -4.2, -0.023)
....				

Here the rows are ordered so that C1 changes most rapidly, followed by C2 and then by C3 so that the grid of cells in the three dimensional C1, C2, C3 space is traversed in a regular order. We can replace this table by one in which only the values for C4 and C5 are included explicitly. The information about the binned C1, C2, and C3 datasets are stored in the descriptions of the structure of quantities C4 and C5. C4 is a normal 3-dimensional image; the pixels of the 3-dimensional array of values in the C4 column are mapped to values of C1, C2 and C3, which are called the axes of the image. C5 is a more complicated object, an image whose pixels are vector-valued. Support for objects like C5 (arrays of vectors) is new, but gives added consistency to the data model. Arrays of vectors are useful, for instance, when the varying centroid position of a source is measured as a function of several parameters.

1.10 Descriptors and Elements and Components

The building blocks of our data are called Descriptors and Elements. The Descriptor represents a named quantity which has an array of Elements associated with it; the simplest case is when the array is trivial and there is only

one Element for the Descriptor. The Element is the value of the Descriptor, sometimes a simple scalar value but in general itself a vector. Each member of the vector has its own name - for instance a descriptor EQPOS representing the equatorial position of something, with a dimensionality of two, has components called RA and Dec. In the original DM design these components only existed as names, and the data they pointed to was accessed through the EQPOS descriptor. In the final implementation, practical convenience drove us to create implicit ‘component descriptors’ for each member of a vector descriptor, so that RA and Dec themselves are scalar descriptors. In general, any vector descriptor will have associated scalar component descriptors for each of its components.

Let’s consider a simple physical quantity: the energy of the Fe K line, which we wish to store as an object FE_K_ENERGY. Suppose we have measured it to be $6.4 \pm 0.3 \text{ keV}$. We will store the name FE_K_ENERGY and the unit keV as part of a Descriptor of real data type. Associated with this Descriptor is an Element of dimension 3: the values 6.50, 5.5 and 7.3 representing the main value and the uncertainty range. We store the range since this lets us easily handle the case of upper limits: an upper limit is just an element for which the lower bound of the uncertainty range is zero or negative. If we get a whole set of measurements of FE_K_ENERGY, we retain the single Descriptor and associate many sets of values with it. In the current release, the DM has no knowledge of the semantics - that the range is actually an uncertainty on the value - and cannot support the case where there is a whole table of values sharing a single uncertainty range in the header. The new design will support the concept of ‘element type’ which encodes this knowledge, although it may be a while before the software is upgraded to do anything useful with the information.

Another type of Descriptor is a Filter descriptor, which has an range but no value. Filter descriptors are used to describe filters, regions, intervals, etc.

Multiple values associated with a single Descriptor are called Arrays. Arrays of scalar Elements are familiar; arrays of vector Elements are more complicated, but are sometimes needed. The simplest kind of array is a one-dimensional array, which simply has a given number of Elements. Note the difference between an array with dimensionality 1 and dimension n (n 1-dimensional Elements), and a vector with dimensionality n and dimension 1 (1 n-dimensional Element).

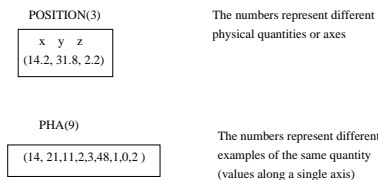


Figure 2: Difference between a vector and a 1-D array. In the first case, each component has a name (e.g. ‘y’); you would plot the n-tuple as a single point in n-dimensional space. In the second case, the different components do not have names. You would plot this as 9 different points along a 1-dimensional space. We also use arrays of vectors: for example, PSF centroid position versus energy and off axis angle.

1.11 Stacks

To work more effectively with multiple sets of data, we introduce the concept of stacks. The simplest stack is just a list of files. However, a more powerful kind of stack is a table one of whose columns contains filenames: in other words, we have a list of files which is **labelled** by the other columns. As an example, let us consider a set of point spread function calibration images which have been taken at some quasi-random set of energies and off axis angles and have similarly random filenames PSF42, PSF13, PSFA1, etc. We make a table PSFSTK as follows:

```
ENERGY  THETA  PSF_FILE
real    real   file
```

```

0.3    42.1    PSF42
0.3     0.1    PSFA1
....
5.2     0.2    PSF13

```

This gives us a 'library' of PSFs which we can look up as a function of the two parameters ENERGY and THETA. If the ENERGY and THETA parameters are table attributes (header keywords) in the individual PSF files, we can imagine a program which would make this stack file PSFSTK automatically by saying: look at all the files in this directory, and for each file with a table attribute OBJECT whose value is equal to 'PSF', add a record to the stack labelled with the values of the table attributes ENERGY and THETA. I will call this operation 'stacking (a set of tables) on ENERGY and THETA'.

We then define a 'stack operation' at the tool level as follows: if the effect of a tool T on a non-stack file F is to make a multi-line table T(F), then the effect of the tool on a stack is to make a new stack table where each entry F in the stack column is replaced by the name of T(F). If the effect of T is to make an output file with a single line, then the entry F is replaced by the contents of that line (so the output file is no longer a stack but a single table).

To continue the earlier example, consider two tools T1 and T2, where T1 takes the histogram of the image pixel values, and T2 returns a one-line table containing the centroid position and total counts. Running T1 on PSF42 makes a new file PSF42_IMHIST (say) with several rows and columns. Running T2 on PSF42 makes a new file PSF42_CTR with several columns but only one row:

```

XCEN  YCEN  TOT_CNTS
real  real  integer

42.3  121.2  141412

```

Then running T1 on PSFSTK should make a new stack as follows:

```

ENERGY  THETA  IMHIST_FILE
real    real   file

0.3     42.1   PSF42_IMHIST
0.3     0.1   PSFA1_IMHIST
....
5.2     0.2   PSF13_IMHIST

```

as well as making all of the individual IMHIST files. But running T2 on PSFSTK should make a single file

```

ENERGY  THETA  XCEN    YCEN    TOT_CNTS
real    real   real    real    integer

0.3     42.1   42.3    121.2   141412
0.3     0.1   52.1    1109.1  32821
....
5.2     0.2   9212.2  104.2   1821

```

The power of this is that it allows us to do aggregate analysis easily: we can now use the generic plot tool to plot, say, XCEN versus THETA to see how those two parameters vary with each other.

This enhanced stack capability will be implemented in the new DM design via special tools, rather than being part of the intrinsic DM syntax.

2 Some general requirements

2.1 Data Model and files

We require that the data model reflect the structure of our science data as generally as possible. Our paradigm for analysing data involves applying tools (programs) to one or several input data files, and generating output files. Data files may be ‘standard data products’ whose structure and contents are predefined in detail by the CXC, ‘user-derived data files’ which follow our general paradigm but whose detailed structure is specified by the user, and ‘compatible data files’ which are produced by external analysis systems (including archives of older missions) but which are sufficiently similar in structure that our software can recognize them. It turns out that almost all our data can be described in terms of instantiations of a single kind of object, which I will call a DM block (or DM Table). There is also a special flavor of DM block called a DM Image which is treated separately in some cases.

A requirement is that the division of our data into separate files should ‘make sense’ to the scientist, logically related information being kept together. An obvious way to do this in the object-oriented paradigm is that each file should contain exactly one DM Table. However, this isn’t the way that files are made by many other software systems, so we have to support a more general approach.

We require that the data model allow the applications programmer to ignore the details of the specific file format conventions (e.g. FITS, QPOE) but also allow some measure of override access to the specific file format writing kernels. At least three kernels will be supported by the model, to support writing and reading ASCII text files, FITS files and IRAF files. By IRAF files I mean IMH files and PROS QPOE files.

2.2 Compatibility Requirements on FITS kernel

We require that as many of the following existing archival FITS datasets should be readable by the FITS kernel as valid Science Datasets: Event lists and XSPEC-type PHA and response matrix files for the following missions: Einstein, ROSAT, ASCA, XMM-Newton and XTE. This imposes requirements on the FITS keywords used to map data model structures.

2.3 The native data model in FITS

FITS files contain a set of independent Header Data Units (HDUs). There are several flavors of HDU but the most important ones are IMAGE and BINTABLE. We will consider a FITS file containing only IMAGE and BINTABLE HDUs. The HDU consists of a header and a data section.

- The header consists of an arbitrary number of header cards.
- A header card contains a keyword (an 8 character case-insensitive string), a value (of one of a number of data types), a comment (a string which is usually ignored by software), a data type (which is not given explicitly, but may be deduced from the formatting of the value), and a card type (deduced from the keyword name).
- The card types are: Mandatory cards, standard reserved cards, local reserved cards, and ordinary cards. By local reserved cards I mean cards whose keywords are not reserved from the point of view of the FITS standard but which are given a reserved meaning in some extra convention to which the particular FITS file adheres. An example is the WCS convention which has been proposed for inclusion in the standard and reserves the meaning of several extra keywords.
- An IMAGE data section consists of an n-dimensional array of numerical values. Associated with the IMAGE data is the data type of the array elements, the number of dimensions, and the size of each axis. This information is contained in reserved header cards; scaling and unit information about the data and coordinate information about the axes may also be associated with the image in this way.

- A BINTABLE data section consists of a set of columns. Each column has a data type and a name, and possibly a unit and various coordinate information. All the columns have the same number of entries.

2.4 Interaction of Data Model and other infrastructure

The data model affects the other infrastructure parts as follows:

- The filtering language describes a virtual dataset in terms of a preexisting one. This description should be complete in the sense that it fills in all the interface requirements of the data model.
- The data model I/O works on filtered (virtual) datasets.
- The data model is pretty much decoupled from work like the scripts, Sherpa parsing, etc.

3 The CXC Data Model, SDS Version 2.0

In this section I present the detailed design for the data model, developed from the earlier model which was based on extensive discussions with Martin Elvis, David Van Stone and Peter Patsis. The model represents a very general kind of table, whose columns can contain vectors or multidimensional arrays, with associated coordinate systems and other metadata. Further metadata can be associated with the table as a whole or with individual columns, and a 'data subspace' indicates the range of values for which the table is valid.

This version 2.0 (Oct 2001) has been significantly updated from the 1994-1997 versions which describe the design prior to implementation of the first DM release.

A note of explanation is required for object-oriented fans (others may skip this paragraph). A Rumbaugh diagram for our design, shown below, indicates that only a small number of distinct objects are used. However, I feel that the Rumbaugh methodology, at least as I have been made to understand it, obscures understanding of the true structure of our data in which multiplicity and aggregation of instances play a key role. I therefore use a slightly different kind of diagram, which I will call a structure diagram, which includes structural components which are not necessarily distinct objects in the OO sense, and which shows separate instances of an object if (and only if) the object is instantiated in a separate role. All of the associates represent 'has a' relationships.

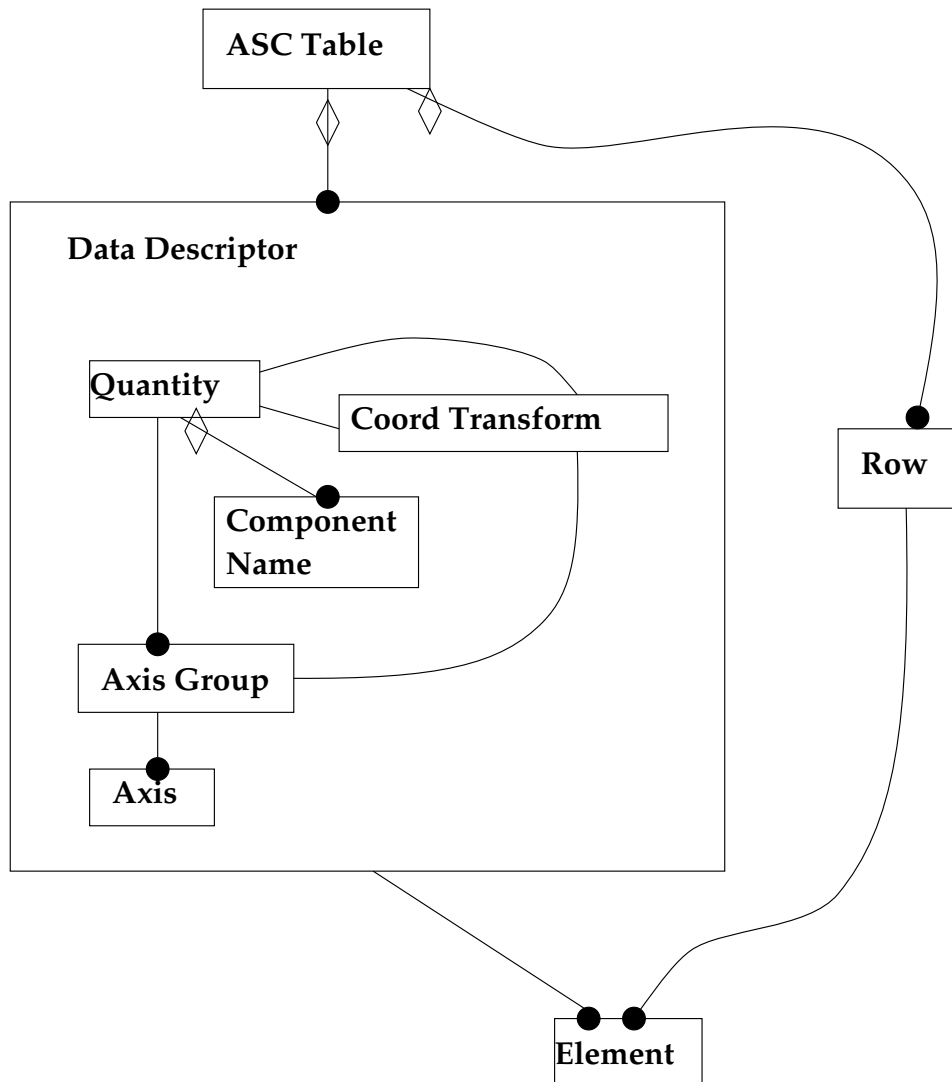


Figure 3: Rumbaugh diagram for ASC Table/block Data Model, showing the fundamental different object classes. The Data Descriptor class is a composite whose components are shown within an enclosing box.

In less technical language, I'm trying to present an abstraction of a particular kind of scientific dataset. The diagrams I show are an attempt to illustrate the different components that go to make up this abstraction. Each box is one of these components, and a line going out of the bottom of one box into the top of another indicates that the second box is a component of the first. A symbol by the end of the line indicates the number of such components that will exist. For instance, in the Table Model diagram the symbol c appears next to both the line from Table Data to Column Data Descriptor and the line from Table Row to Column Data Cell. This indicates that there are the same number of Column Data Cells in a Table Row as there are Column Data Descriptors in a Table Data, and that we will represent this number as c (it happens to be the number of columns in the table). If no multiplicity is indicated on a line, there is exactly one of the components attached to the parent object.

Text underneath a horizontal line below the name of the object indicates parameters (attributes) associated with that object. For instance, the DM Table has a Name. A double line at the bottom of an object box indicates that the structure of the object is complex and there's a whole separate diagram for it later on. For instance, DSS Descriptor is a particular kind of Descriptor object; the Descriptor has a diagram of its own, and the text below defines what kind of a Descriptor a DSS Descriptor is (For instance, unlike a general Descriptor, it can't be an Array.)

All multiple subcomponents are considered to be ordered. In other words, there is a defined order of the columns in a table, a defined order of header keywords in the header, and a defined order of the axes in an array. One may refer to a column by its name (e.g. TIME) or by its number (e.g. Column 4). A single Table Column may map to several columns in an underlying table format (e.g. FITS BINTABLE), and in general the numbering of a Table component is distinct from the numbering of the corresponding structure in the underlying data format.

3.1 DM Table

The highest level object is the DM block. It consists of three main parts: the Table Data proper, the Header, and the Data Subspace. Each of these contains Data Cells made up of arrays of Components which contain the actual data, and Data Descriptors which provide metadata about the meaning of the Components.

The DM block has a single attribute of its own: the table name.

4 Table Data Section

4.1 Table Data

The Table Data Section represents a table with r rows and c columns. The intersection of a row and a column is a Data Cell; all the Data Cells in a column have the same structure, and contain the same type of data; they are described by the Data Descriptor for that column. The different columns in a row may have different structures.

Associated with the Table Data section is an ordered list of Preferred Columns, as a hint to generic software which only operates on a given number of columns without specifying specific column names.

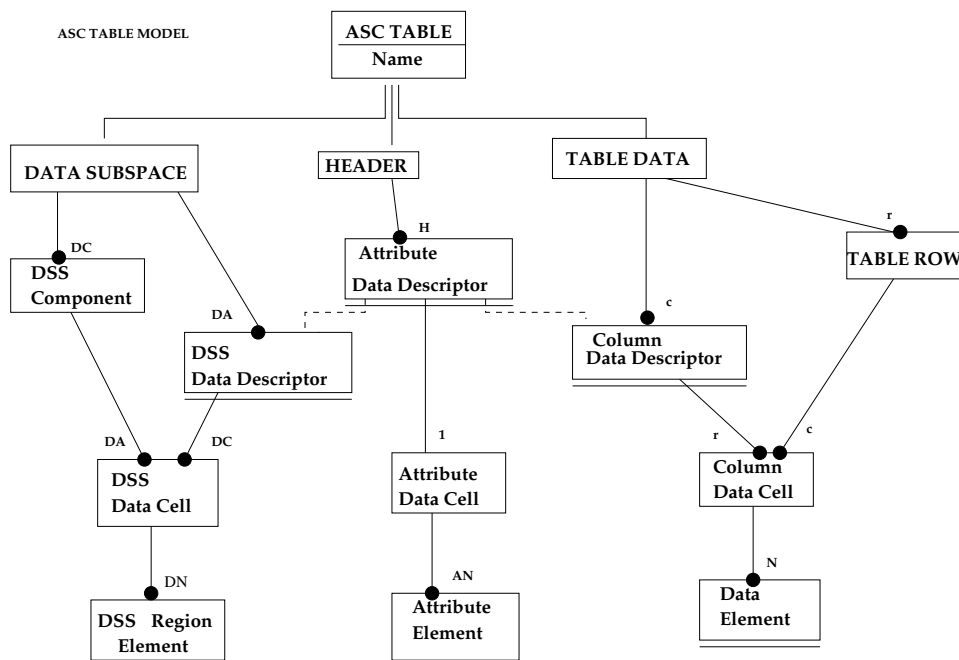


Figure 4: Data Model 1: Overall structure of the data model, showing the ASC Table, used as the highest level data object encapsulating all others.

4.2 Descriptor

The Descriptor object is proved to describe the structure and properties of named quantities. It is basically a structure which provides a name, a unit and other descriptive information, and specifies a data type. It is the abstraction of the FITS BINTABLE's TTYPE*n*, TUNIT*n*, etc., header keywords for a column.

The Descriptor has:

- Name
- Unit
- Description
- Data type
- Display format
- Element dimensionality d
- Element type t
- Array dimensionality n
- Array specification: dimensions $n_1 \dots n_n$ grouped into Axis Groups, total number of values $N = \text{product of the } n_i$.
- Component descriptors c_1, \dots, c_d
- Coordinate descriptors C_1, \dots, C_m
- Axis Group descriptors
- Parent descriptor
- Elements $E(n_1, \dots, n_N)$
- Kernel marker

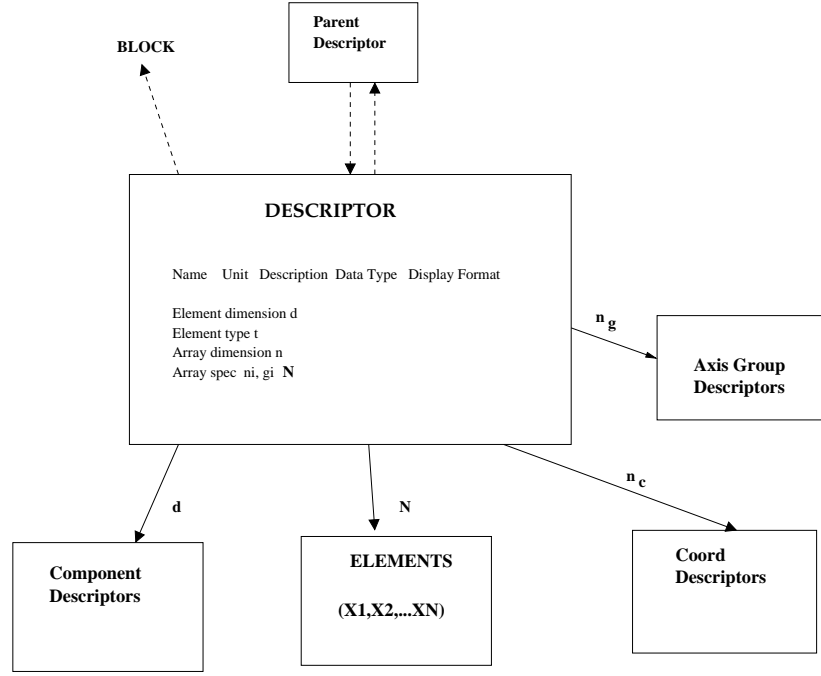


Figure 5: Data Model 2: The Descriptor object, with attributes and pointers to associated Descriptors and Elements.

In more detail, the attributes are:

- The descriptor name, a character string. Any ASCII character string shall be supported, but we recommend that the string shall consist of only alphabetic upper or lower case letters a-z,A-Z; numeric digits, 0-9; the symbols +,- and underscore (_). In particular, spaces are not permitted (except trailing spaces which are not considered to be significant). Internally defined descriptors in the DM use the leading character #. At this level, case is significant, although we anticipate that user access routines will not be case-sensitive and recommend that names be unique within a table even when case is ignored. The idea here is that we may want to name something MaxVoltage instead of MAXVOLTAGE so that the software knows how to print it nicely, but we don't want to require that the user has to get the capitalization right when searching for it. So case is remembered, and returned correctly, but matches are case insensitive.
- The unit. A character string which specifies the physical unit. Should comply with the HEASARC/OGIP format on unit strings or the JCMLIB specification for unit strings.
- The Description is a string which is used to label human readable output such as ASCII print files and graphical axis labels. It is a longer name which may include spaces and other special characters, including backslash. I suggest the use of TeX escape sequences which are supported by some graphics libraries such as SM, for instance '\ alpha' for α . Thus a Descriptor might have the name 'RA' and the description '\ alpha (J2000.0)'.
I suggest the use of TeX escape sequences which are supported by some graphics libraries such as SM, for instance '\ alpha' for α . Thus a Descriptor might have the name 'RA' and the description '\ alpha (J2000.0)'.
- The Display Format indicates the preferred output format for a single data value associated with the quantity in a text browser. It is required that the Display Format can be returned as a string Fortran format specification compatible with the TDISPn keyword in a FITS file. This optional information may be provided as a hint to browsers to let them format tabular output efficiently. For instance, a quantity stored as a 4 byte integer might be known to only take values less than 1000, allowing a display format of 'I4' instead of the larger 'I10' needed by an arbitrary 4 byte integer. Pixel coordinates might be displayed as 'F8.3' while a time specification in seconds might require the greater precision of 'F20.6'. However, the display format may be absent or the browser may choose not to use it, it's just provided to help make the output pretty.
- The Data Type indicates the type of data in an associated Element. Supported data types shall include:
 - Integer, 2 byte
 - Integer, 4 byte
 - IEEE Real, 4 byte. The specification of IEEE here indicates that it must be possible to return the data in IEEE format, and it must be possible to store IEEE special values such as NaN and -Inf. How the data are actually stored internally or in a data file is an implementation detail.
 - IEEE Real, 8 byte
 - Logical, 1 byte
 - String, specified fixed number of ASCII bytes s (Unicode is not currently supported).
 - Bit, 1 byte
 - Unsigned Byte, 1 byte
 - Unsigned Integer, 2 bytes
 - Unsigned Integer, 4 bytes

In addition, the following data types are under consideration for future support:

- Extended Unsigned Integer, 8 bytes
- Complex, 16 bytes

Table 1: Codes for Data Types

Data Type	API routine suffix	FITS CFORM	FITS TFORM
Integer/2	s	'I'	'I'
Integer/4	l	'J'	'J'
Real/4	f	'E'	'E'
Real/8	d	'D'	'D'
Logical	q	'L'	'L'
String	c	'A'	'A'
Unsigned/1	b	'B'	'B'
Unsigned/2	us	'U'	'I'
Unsigned/4	ul	'V'	'J'

- The Element Dimensionality specifies the dimensionality d of all Elements associated with this Descriptor. The default is $d = 1$.
- The Element Type can be Value (V); Value with Uncertainty (U), Value with Fixed Uncertainty (UF), 2D Region (REG), Bin (BF), Bin Start (SF), etc. The different element types are discussed in full in the section on elements. All Elements associated with the Descriptor must be of the same Element Type. Element types are newly supported in DM2.0.
- A Kernel marker. This is a placeholder to support extra information needed to reconstitute a clean file for a particular kernel.
- Array Dimensionality n specifies the dimensionality of the array of Elements making up a single Cell associated with the Descriptor. If $n > 0$, there must be an Array Specification associated with the Descriptor; if $n = 0$ there is no Array Specification. All Cells associated with the Descriptor must have the same array dimensionality and array specification.
- If $d > 1$, there are a set of d Components. For instance, we might define a 2-dimensional Descriptor with name SKYPOS and component names RA and Dec. If $d = 1$ then the single component is defined to be identical with the Descriptor.

We name certain special cases:

- – A Descriptor with $d = 1$ and $n = 0$ is known as a Scalar Descriptor.
- – A Descriptor with $d > 1$ and $n = 0$ is a Vector Descriptor.
- – A Descriptor with $d = 1$ and $n > 0$ is a Scalar Array Descriptor.
- – A Descriptor with $d > 1$ and $n > 0$ is a Vector Array Descriptor.
- Finally, a descriptor may have a Comment field. This Comment field consists of arbitrarily many 72-byte text strings each with an associated 8-byte tag. The default value of the tag is the string 'COMMENT'. Other values of the tag are not guaranteed to produce valid files for all kernels, although 'HISTORY' and blank are valid for FITS files. The Comment field text may appear in the underlying file header anywhere following the appearance of the descriptor name and preceding the next descriptor name.

4.3 Array Dimensions

An Array Dimensions specification describes the arrangement of a set of N elements into an n -dimensional array. The n axes of this array, $i = 1, \dots, n$, have dimension (size) n_i , so that

$$\prod_{i=1, n} n_i = N$$

The elements $E(p_1, \dots, p_n)$ of the array are labelled by array pixel numbers, which are an ordered n -tuple $P = (p_1, \dots, p_n)$.

4.4 Array Axis

Each axis i of the array is defined by a given dimension (size, number of pixels) n_i . We adopt the FITS (and Fortran) convention in which the pixel numbers start at one, and in which a default storage order is implied in the following sense: an Element Number e is defined equal to

$$e(P) = p_n + (p_{n-1} - 1) * n_n + ((p_{n-2} - 1) * n_n n_{n-1} + \dots$$

or

$$e(P) = \sum_{i=1}^n \left((p_i - 1) \prod_{j=i+1}^n n_j \right)$$

where $\prod_{j=n+1}^n n_j$ in the final term of the sum is interpreted to be equal to one. A mechanism will be supplied to return the elements in element number order. In FITS files and in Fortran arrays, the array elements must be actually stored in element number order.

4.5 Axis Groups

We add a little extra structure to the array to group axes which may have common coordinate transforms. Axis Groups are to image axes as Vector Descriptors are to table columns. In our model we consider something like detector pixel position to be a single, two-dimensional, Descriptor; if we have a data cube of detector pixel position DETX, DETY versus energy E we wish to emphasize the fact that DETX and DETY are related to each other in a way that they are not related to E. In this view, the three dimensional data cube DETX,DETY,E is instead a two dimensional array with two axes DETPOS and E, in which the first axis is itself two-dimensional. This first, two-dimensional axis may have a coordinate system on it which applies a two-dimensional spherical rotation, or it may have a mask on it which specifies a two-dimensional region, in each case requiring that treatment of DETX and DETY be coupled. In contrast, we do not expect to get situations where we must treat DETX and E in a coupled way (if we do, they will have to be treated at a higher level).

The Array Dimensions specification adds the concept of Axis Groups. In the example above, the three dimensional array has two axis groups, one a two-dimensional axis group containing the first two axes and another one-dimensional axis group containing the third. We can label the array by axis group pixel numbers $P_G = ((p_1, p_2), p_3)$, an ordered pair of a two-dimensional detector position pixel and an energy bin.

Say there are g axis groups each of dimensionality $g_m, m = 1, \dots, g$. We have

$$\sum_m g_m = n.$$

4.6 Parent Descriptor

If the Descriptor is a component Descriptor, an Axis Group or Coordinate Descriptor, etc., it has a parent - another descriptor which refers to it. A raw table column may have no parent descriptor; filtered table columns point back to the raw column.

4.7 Coordinate Transform Descriptor

The Elements of a Coordinate Descriptor are defined implicitly by mapping the Elements associated with another Descriptor using a Coordinate Transform. A simple example is the mapping of mission time TIME in seconds to Julian Date JD in days. We define this transformation by choosing a reference value of TIME (usually 0.0) and the corresponding reference value of JD (the JD when TIME is equal to 0.0; say 2445200.0), and defining the transformation relative to this reference value. If TIME is the correct time in seconds since the reference value, then the transformation type is LINEAR and the transformation scale is 1.0/86400.0 (the number of days in a second); this completely determines the transformation. If TIME is a spacecraft clock with glitches and resets, the transformation may be a lookup table or a polynomial with a more complicated definition. Lookup transforms are not yet implemented in the DM.

In general, we consider a coordinate transform to link two Descriptors which have the same Element Dimensionality d . One Descriptor is referred to as the Parent or Pixel Descriptor and one as the Coordinate Descriptor (this does not necessarily imply that the Pixel Descriptor has units of pixels; the names evoke the FITS keywords CRPIX and CRVAL). The transform consists of a Coordinate Transform Specification which has a Transform Type, a set of d Transform Scales Δ_i ($i = 1, ..d$), and associated Transform Parameters specific to the transform type. It also has a Reference Pixel Element and a Reference Coordinate Element, which are Value Elements (Elements of type V, see below) corresponding to the Pixel and Coordinate descriptors. In our example above, the Pixel and Coordinate Descriptors are TIME and JD, and the Pixel and Coordinate Elements have values 0.0 and 2445200.0.

The point here is that we choose to represent an arbitrary transformation by a local linear transform about a reference point, plus higher order corrections. This has three advantages: it maps directly to the FITS CRPIX/CRVAL/CDELTA convention; it ensures that we have a defined 'center' for our transformation, which can be used as a default location by an application; and often the transformations we use are linear, and don't require any higher order parameters, so it makes the usual cases simple.

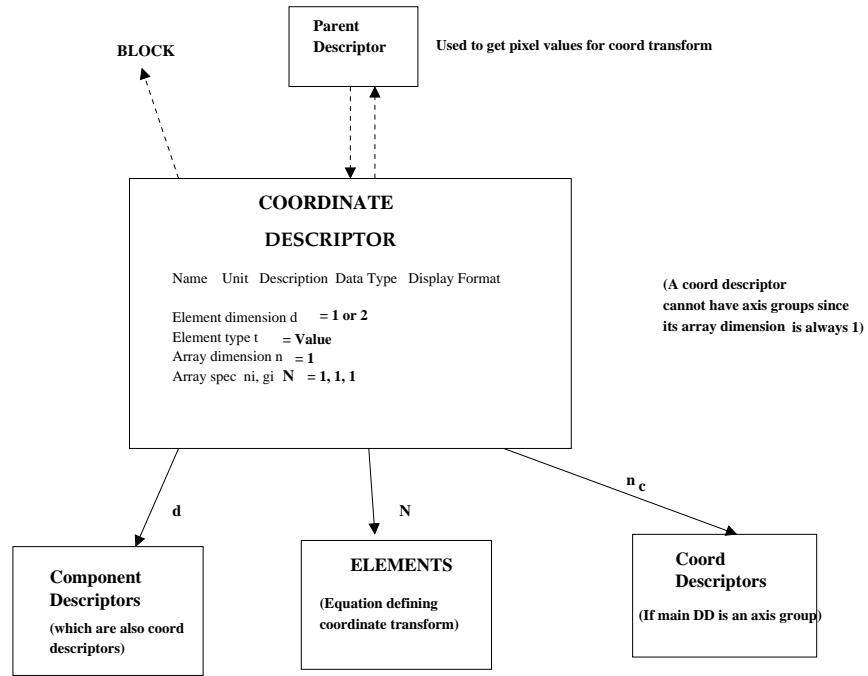


Figure 6: Data Model 3: Use of the Coordinate Transform

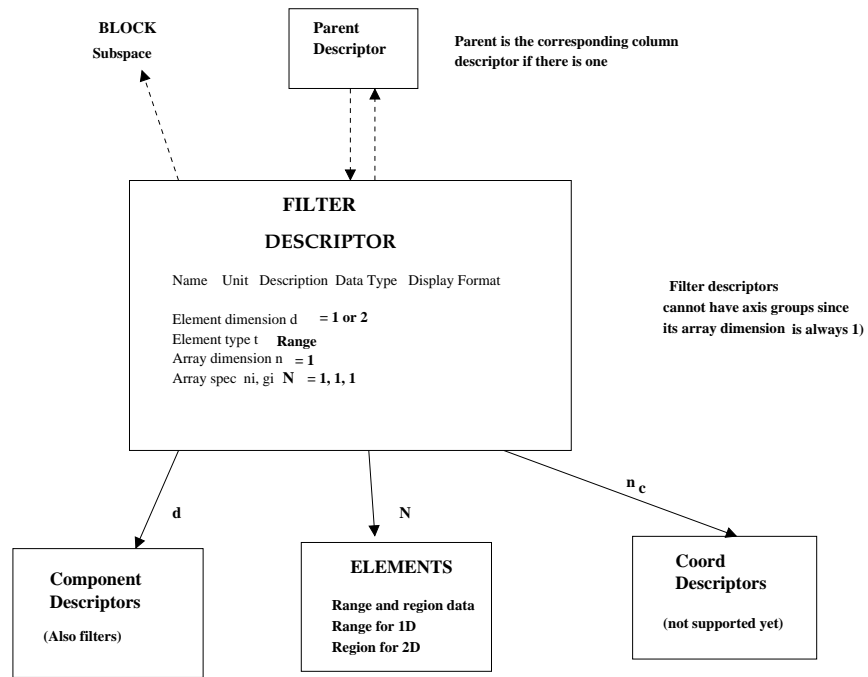


Figure 7: Data Model 3b: Filter Descriptors

4.8 Column Data Descriptor

A Data Descriptor provides information about a Descriptor which we're going to provide values for. The simplest, minimal Data Descriptor is a Data Descriptor which is a scalar Descriptor. More complicated Data Descriptors provide support for vector Descriptors, for Arrays (a Descriptor with an Array Specification), and for associated coordinate and axis descriptors.

Every Data Descriptor has a single Data Descriptor. A Data Descriptor whose Data Descriptor is a Vector Array Descriptor is called a Vector Array Descriptor, and so on. If the Data Descriptor is an Array Descriptor ($n > 0$) then there is an associated Array Specification with Axis Groups and Axes. A Scalar Data Descriptor in a Column Data Descriptor is the abstraction of the FITS keywords TTYPE n , TFORM n , TUNIT n , etc. An Array Data Descriptor corresponds to the FITS BINTABLE multidimensional array convention for TFORM values. Vector Data Descriptors do not correspond to any existing convention in FITS.

Associated with the Data Descriptor there may be a Data Coordinate Descriptor linked to it by a Data Coordinate Transform. For instance, a table may have a column TIME with values included explicitly in the table cells. The TIME column may have associated with it a Descriptor JD which gives the Julian Date. The individual values of JD are not stored explicitly, but are implied by the JD to TIME coordinate transform. JD is a Data Coordinate Descriptor associated with the Data Descriptor TIME. The Data Coordinate Descriptor and Transform are the abstractions of the FITS keywords TCTYP n and TCRVL n , TCRPX n , etc.

A Column Data Descriptor with a Data Descriptor which is an Array has an Array Specification with one or more Axis Groups. Each Axis Group may have an associated Axis Group Descriptor, related to it by a coordinate transform called a Pixel Coordinate Transform which must be of transform type LINEAR. The Axis Group Descriptors are the labels of the axes of the array. For instance, we may have a Data Descriptor PSF which is a three dimensional array with axis groups $g_1 = 2$ and $g_2 = 1$, associated with Axis Group Descriptors DETPOS ($d = 2$, component names DETX and DETY) and ENERGY ($d = 1$). The element dimensionality of the Axis Group Descriptor must be the same as the dimensionality of the Axis Group.

Further, the Axis Group Descriptors may themselves have associated Axis Group Coordinate Descriptors related to them by Axis Group Coordinate Transforms. Consider another example in which the Array has $n = 2$, $g = 1$, $g_1 = 2$ and the single Axis Group Descriptor is SKYPOS with components X and Y representing the X,Y sky pixel coordinate positions. We may associate with it an Axis Group Coordinate Descriptor EQPOS with components RA and DEC, linked by a coordinate transform of type TAN, representing the actual equatorial sky positions. The Axis Group Coordinate Descriptors are the abstractions of CTYPE n in a FITS image, while the lack of support for Axis Group Quantities themselves (such as SKYPOS X,Y) is an unfortunate limitation of current FITS practice.

4.9 Interval type

An Interval defines a contiguous subset of the data values of the appropriate data type. Intervals are only meaningful for data types where a well defined ordering of the data values exists. For integer and real types this is the usual ordering; for string types this is defined to be the ASCII ordering.

The most general Interval is a minimum value, a maximum value, and an interval type. Possible interval types are closed, open, semi-open lower, and semi-open upper, denoted as [a:b], (a:b), (a:b], and [a:b) respectively. These are defined as:

$$\begin{aligned}x \in [a : b] &\Leftrightarrow a \leq x \leq b \\x \in (a : b) &\Leftrightarrow a < x < b \\x \in (a : b] &\Leftrightarrow a < x \leq b \\x \in [a : b) &\Leftrightarrow a \leq x < b\end{aligned}$$

The semi-open intervals are useful for ensuring that boundary values are not counted twice. For integer and string data types, the only possible type of interval is Closed. This is also the default interval type for real data types.

In the DM2.0 design, Intervals are assumed to be [a:b) in all real cases and [a:b] in all integer cases. Adding support for explicit control of interval type is under consideration.

4.10 Elements

The actual data for the table is stored in Elements. An Element must be associated with a Descriptor. A single Element contains values for one instance of the Descriptor. For example, if the Descriptor TIME has element type Value with Uncertainty (VU) and element dimensionality 2, with components TIME and TIME_UNC, then an Element associated with TIME has one value of the TIME and one value for TIME_UNC. If the Descriptor DET has element type Value (V) and element dimensionality 2, then a single Element of DET has two Value Elements. The simplest kind of element is an element of type Value and dimensionality 1, which is a single value (numeric or string according to the associated Descriptor's data type.)

The special element type REG applies only to 2-dimensional elements and is a string defining a region in PROS Regions syntax. With the exception of this element type, all d-dimensional elements consist of uncoupled element components for each of the d dimensions. The most general element component is a Value plus its Uncertainties or Ranges.

Eventually, we propose to support three different uncertainties: statistical, systematic zero point, and systematic scale. In addition, we define a total uncertainty which is a function of these three. We also use the same paradigm to record bin ranges. Our approach is to treat the systematic uncertainties as separate add-ons, with our default description being a single value and uncertainty, which is to be interpreted as a statistical uncertainty if the systematics are present and as a total uncertainty if they are not.

- If no uncertainty at all is present, the code is V (Value).
- The most flexible representation is the Interval Uncertainty (I) which uses an Interval to define the minimum and maximum values within the significance range. If the minimum value is zero or less, the measurement is termed an upper limit. If the Value component has value v , and the Interval has min and max of v_1 and v_2 , then for a closed interval type

$$v_1 \leq v \leq v_2.$$

Note that the range center $(v_1 + v_2)/2$ is not necessarily equal to v .

- A second, more common representation is the Two Sided Uncertainty (T), in which the offsets σ_+ , σ_- from the nominal value are given. This has the advantage that it may be often used as a Fixed Uncertainty. In terms of the Interval Uncertainty,

$$v_1 = v - \sigma_-, v_2 = v + \sigma_+.$$

- The One Sided Uncertainty (U) is the same as the two sided, but both upper and lower uncertainties are equal.

$$v_1 = v - \sigma, v_2 = v + \sigma.$$

- The Bin (B) is the same as the one sided uncertainty, but the full bin width w rather than the half bin width σ is given. This is more usually employed when the interpretation is a binned dataset rather than an uncertainty.

$$v_1 = v - w/2, v_2 = v + w/2.$$

- The Bin Start (S) is the same as the Bin, but the Value is deemed to be the start of the bin rather than the center:

$$v_1 = v, v_2 = v + w.$$

This representation is often used for light curves.

- The Range (R) is the same as the Interval Uncertainty but there is no associated Value. If a Value is required, it is assumed to be $v = (v_1 + v_2)/2$.
- The scale uncertainty is always represented as a single nonnegative dimensionless real value (K) so that the implied range around the value v is

$$v_1 = v(1 - K), v_2 = v(1 + K).$$

- We also want to support a two sided scale (L) with different upper and lower scale errors, which arises when we take the logarithm of a Descriptor with different upper and lower uncertainties.
- Finally, sometimes data is just provided in the form of detections and upper limits. We define an element type Z which consists of a value v_d and a limit flag f , with the meaning

$$\text{if } f \text{ then } v = v_d \text{ else } v_1 = 0, v_2 = v_d.$$

However, I don't propose that we support this element type initially.

Each of these range types can be Fixed, in which case we append the letter F to the element type. We will further require that elements in a Table Column have a fixed Interval Type for all cells of the column.

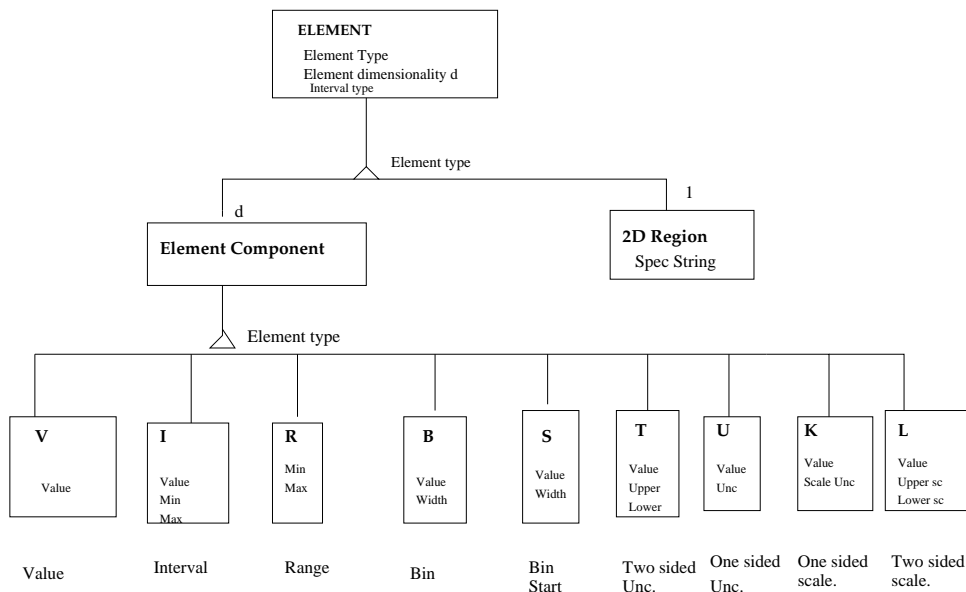


Figure 8: Data Model 5: The Element object, used to store the actual values. There may be many elements described by a single Descriptor. There is one Element component for each dimension of the element, except if systematic uncertainties are included in which case there may be up to three Element components for each dimension.

We will later add to the Descriptor object a systematic zero point uncertainty type and a systematic scale uncertainty type, the default values of which are null (not present). The legal values are the same as for the Element type, and if they are present the usual values are UF for the zero point uncertainty and KF for the scale uncertainty.

4.11 Region Description

For the two dimensional region descriptions we would like to support those in current systems, namely:

- Bitmap: appropriate for a binned dataset, provides a list of the pixels in the region.
- Polygon: an ordered list of n points describing a closed polygonal region.
- Shape: A parameterized shape, including the cases Circle, Annulus, Ellipse, Box, Pie.

Bitmaps are not yet supported.

We can describe the Shape with the following parameters:

- Shape type: elliptical or rectangular.
- Shape center x_0, y_0 .
- Shape radial range r_1, r_2 , interval type. If $r_1=0$, have a Circle or Box. If $r_1>0$, have an Annulus or annular box.
- Aspect ratio a , ratio of major to minor axis. If $a=1$ have a circle or square; if $a<1$ have an ellipse or rectangle.
- Shape orientation θ_0 ; measures angle between major axis and x axis. Irrelevant if $a=1$.
- Shape azimuthal range θ_1, θ_2 , interval type. The default is $\theta_1=0$ and $\theta_2=360$ deg. Any other value gives you a pie or sector (for shape type elliptical; shape type rectangular may not support sectors).

4.12 Table Data Cell

A Data Cell is associated with a Data Descriptor and contains one set of Elements for that Data Descriptor. The number of elements in the cell is equal to the number of elements in the array specification for the data descriptor; in particular, if there is no array specification (data Descriptor array dimensionality equal to zero) there is exactly one element in the cell. The elements in the cell can be accessed via pixel number or element number as discussed in the section on array specifications and axes.

4.13 Table Row

In a Table Data section, there is some specified number r of Table Rows. Each Row may be thought of as containing one Data Cell for each of the Column Data Descriptors. More precisely, there is one Data Cell associated with each combination of row and column.

4.14 Descriptor Groups

A descriptor group is a simple object with a name and an array of descriptor pointers. It allows users to manipulate and refer to related collections of descriptors. API and kernel routines will be provided to define and access such descriptor groups.

5 Data Subspace

5.1 Introduction

What distinguishes a photon event list from a table in an ordinary database? The rows of the event list represent individual, asynchronous events. They cannot be interpreted without knowing the filter through which those events were selected. Suppose we detect photons only between times 100 and 200. Is this because the source flared during that time, or because the satellite was only looking during that time period? To be more precise, if you just have an ordinary table of rows, what you are missing is the information about what rows would NOT have been allowed in the table - in the photon event list case, which events would NOT have been detected. We are then led to the concept of the data subspace: in the space of all possible data, what subspace is being sampled by the current table?

This idea is closely connected with the idea of filtering. The data subspace is simply the filter that has been applied to the data. However, we're not just talking about user filters applied during processing, but also implicit filters applied by the act of observation at a particular time with a particular instrument. If the user then filters the data further, the new data subspace is simply the intersection of the filter with the old subspace.

If two datasets are merged, the new data subspace is the union of the old ones. In this case, however, we lose some information: the data subspace paradigm doesn't retain information about which of the original subspaces a particular row belonged to. This is the usual problem with binning data together, which we can illustrate with a familiar example: combining two pulse height spectra. Suppose we have two event lists E1 and E2 with the following data, representing events from two different ACIS chips which are distinguished by different ranges of detector position DETPOS:

E1 subspace: DETPOS=[0:1024,0:1024]

E1 table:

DETPOS	PHA	TIME	
100	245	8	4922.2
231	928	17	4812.5
....			

E2 subspace: DETPOS=[1024:2048,0:1024]

E2 table:

DETPOS	PHA	TIME	
1241	621	22	4924.3
1782	212	7	4092.2
...			

If we extract two PHA histograms P1 and P2, retaining only pulse heights from 2 to 100 and selecting a region near the boundary of the chips where we think there is a source, we get:

P1 subspace: DETPOS=[1000:1024,800:825], PHA=[2:100]

P1 table:

PHA	COUNTS
2	0
3	4
....	
100	1

P2 subspace: DETPOS=[1024:1124,800:825], PHA=[2:100]

P2 table:

```

PHA  COUNTS
2    1
3    2
...

```

If we then merge these two datasets to form P3, we get:

```

P3 subspace: DETPOS=[1000:1124,800:825], PHA=[2:100]
P3 table:
PHA  COUNTS
2    1
3    6
....

```

A tool to build the XSPEC response matrix would then check the DETPOS region to see which chips were involved. In the case of P3, it would see that 20 percent of the region was on one chip and 80 percent on the other, and would average the two response matrices in that proportion. We have lost any information about which counts came from which chip. If instead we merge the lists E1 and E2 to form a new event list which retains the DETPOS column, and then filter on position and PHA but don't bin to make the histograms, we get E3:

```

E3 subspace: DETPOS=[1000:1124,800:825], PHA=[2:100]
E3 table:
DETPOS  PHA
1012 814  8
1182 803 18
...

```

although the data subspace is the same as for P3, the information about which chip is involved for a given event is still available via the DETPOS value for the given event.

In general, any tabular data may have a data subspace which describes the range of data for which the table applies. The descriptors in the data subspace are not necessarily the same as the descriptors in the table itself - see the example of P3 above in which DETPOS is in the data subspace but not in the table.

5.2 Unions of subspaces

A more complicated case of merging subspaces is when we wish to use 'incompatible' filters. For example, perhaps the second chip has unreliable data in PHA channels 2 to 10, so we want to apply a different PHA filter to it. We filter E1 and E2 with different filters and then merge them to make E4:

```

E4 subspace:
  Component 1: DETPOS=[1000:1024,800:825], PHA=[2:100]
  Component 2: DETPOS=[1024:1124,800:825], PHA=[11:100]
E4 table:
DETPOS  PHA
1012 814  8
1182 803 18
...

```

When two filters (subspaces) are unioned (logical OR), we describe them as different 'components' of the subspace.

What if the different filters involve filtering on entirely different quantities? Consider the case when E1 is filtered on PHA and E2 is filtered on TIME.

E5 subspace:

Component 1: DETPOS=[1000:1024,800:825], PHA=[2:100]

Component 2: DETPOS=[1024:1124,800:825], TIME=[4823.2:4890.1],[5012.4,5100.0)

E5 table:

DETPOS	PHA	TIME
1012 814	8	4902.54
1182 803	18	4823.80
...		

To simplify the treatment, we note that we can make the quantities involved in the two components the same by adding the trivial filters $TIME=[-\infty : \infty]$ and $PHA=[-\infty : \infty]$ to components 1 and 2 respectively. Doing this lets us store a single list of the descriptors involved in a data subspace, instead of requiring us to maintain separate lists for each component.

5.3 General definition

1. A Data Subspace (DSS) D consists of $DC = 0+$ Data Subspace Components (DSS Components) $C(i), i = 1, DC$ and a list of $DA = 0+$ Data Subspace Data Descriptors or Data Subspace Axis Groups $A(j), j = 1, DA$. (*Note: The notation $n_C = 0+$ means that there are zero or more of the entities in question, and that the number of entities will be denoted by DA .*) There is usually only one DSS Component in a DSS, i.e. $DC=1$. The name Axis Group reflects the fact that the data subspace could be represented by an array with those axis groups (although the pixel values of that array are not defined).
2. A Data Subspace Data Descriptor or Data Subspace Axis Group is a named object which has the same properties as the generic Data Descriptor defined above, particularly including a name and a dimensionality. An example of a data subspace axis group might be TIME, or POSITION. However, a Data Subspace Data Descriptor may not have associated array Axis Group Descriptors, or array Axis Group Coordinate Descriptors. Further, it must have array dimensionality 1. An important distinction between the DD for Table Data and the DD for a Data Subspace is that the array dimension n_1 is to be interpreted as the maximum dimension for any data cell, rather than the actual dimension for each data cell (see below). However, the Data Subspace Data Descriptor is allowed to have a Data Coordinate Transform and a Data Coordinate Descriptor.
3. A Data Subspace Component $C(i)$ consists of DA DSS Data Cells $RV(i, j)$, one for each axis group of the parent data subspace.
4. The Data Cells of a data subspace component consist of $n_R = 0+$ Region Elements $R(i, j, k), k = 1, n_R(i, j)$. An example of such a Data Cell is a set of Good Time Intervals, or a spatial mask consisting of several components. The different Data Cells corresponding to different DSS Components may have different values of n_R , unlike the Data Cells for different rows of a Table Data section which must all have the same array sizes. Since there is usually only one DSS component, this doesn't usually matter.
5. A Data Cell may be defined implicitly as a World Coordinate Data Cell. For instance, if the Data Subspace Axis Group is pixel sky coordinate position SKYPOS (X,Y), and this has a Data Coordinate Descriptor EQPOS (RA,DEC) related to it by a Data Coordinate Transform, then we may express the Data Cell as a set of region elements attached to EQPOS (the Data Coordinate Descriptor) rather than SKYPOS (the Data Descriptor) - say, a circle expressed as '(c 14:04:11 -00:23:12 6.2)', i.e. a 6.2 arcmin circle around the specified sexagesimal RA and Dec, instead of '(c 4212.2 5123.2 42.1)' in pixels. I haven't included this explicitly in the diagrams; in the FITS implementation I have suggested parallel keywords DSn and DSCn for regions expressed in the pixel and world systems respectively.

6. A Region Element $R(i, j, k)$ in a data subspace data cell is a range element if the dimensionality of the corresponding Data Subspace Axis Group is 1, and is a 2D Region Element in the dimensionality of the corresponding Data Subspace Axis Group is 2.

From a set-theory point of view,

$$RV(i, j) = \cup_k R(i, j, k)$$

and

$$C(i) = \cap_j RV(i, j)$$

and

$$D = \cup_i C(i) = \cup_i \left(\cap_j (\cup_k R(i, j, k)) \right)$$

Data Subspace

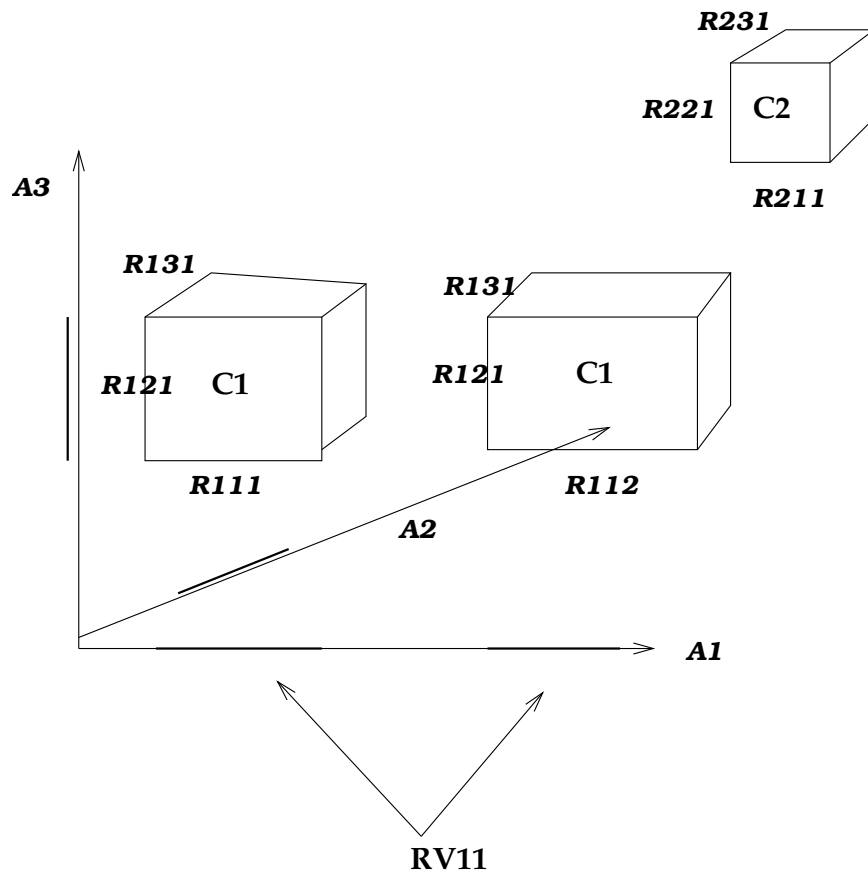


Figure 9: Illustration of a data subspace.

7. A data point P, consisting of values V_j , $j = 1, DA$, is said to be 'in' the data subspace if it is in any one of the components. It is in a component if it is in all of that component's data cells. It is in a data cell if it is in any of the data cell's region elements.
8. The intersection of two data subspaces D^1 and D^2 is calculated as follows: First extend the lists of axis groups of each subspace to be the same. Then

$$D^1 \cap D^2 = \bigcup_i \left(\bigcap_j (\cup_k R^1(i, j, k)) \right) \cap \bigcup_m \left(\bigcap_j (\cup_n R^2(m, j, n)) \right)$$

or

$$D^1 \cap D^2 = \bigcup_i \bigcup_m \left(\bigcap_j (\cup_k \cup_n R^1(i, j, k) R^2(m, j, n)) \right)$$

The case of a single point can be understood as a special case of this. Consider the value components V_j as closed zero-length ranges $[V_j : V_j]$; then P is a data subspace with one component and $R(i, j, k) = [V_j : V_j]$. The above formula tells us to intersect each component with the corresponding range.

Examples of intersection of data subspaces: First, let's take the point case. Let the data subspace be that of E5 above:

Component 1: DETPOS=[1000:1024, 800:825], PHA=[2:100], TIME=[:]

Component 2: DETPOS=[1024:1124, 800:825], PHA=[:], TIME=[4823.2:4890.1], [5012.4, 5100.0)

Then let P be the point (DETPOS, PHA, TIME)=(1100, 812), 200, 5050). We have:

```

A(1) = DETPOS
A(2) = PHA
A(3) = TIME
R(1,1,1) = Box 1000:1024, 800:825
R(1,2,1) = [2:100]
R(1,3,1) = [:]
R(2,1,1) = Box 1024:1124, 800:825
R(2,2,1) = [:]
R(2,3,1) = [4823.2, 4890.1)
R(2,3,2) = [5012.4, 5100.0)
V(1) = (1100, 812)
V(2) = 200
V(3) = 5050

```

So first we intersect P with component 1. The intersection is null, since V(1) has no overlap with R(1,1,1) and V(2) has no overlap with R(1,2,1). Next we intersect with component 2. The intersection of V(1) with R(2,1,1) is V(1) itself; similarly for V(2). V(3) is outside R(2,3,1) but inside R(2,3,2) and thus inside their union as required. So the intersection of P with component 2 of the subspace is P itself. Thus, P is inside the subspace.

Now let's take the intersection of two filters. Let the second space be a simple time filter with two intervals, TIME=[4000:4800],[6000:7000]. To do the intersection we add the missing axes:

```

R(1,1,1)=[:, :]
R(1,2,1)=[:]
R(1,3,1)=[4000:4800]
R(1,3,2)=[6000:7000]

```

Then evaluating the intersection equation gives the expected result:

```
A(1) = DETPOS
A(2) = PHA
A(3) = TIME
R(1,1,1) = Box 1000:1024, 800:825
R(1,2,1) = [2:100]
R(1,3,1) = [4000:4900]
R(1,3,2) = [6000:7000]
R(2,1,1) = Box 1024:1124, 800:825
R(2,2,1) = [:]
R(2,3,1) = [4823.2,4800.0]
```

Note that the second element of the TIME region vector in component 2 has disappeared, since it had no overlap with the new filter. The interval type of the first element has changed, it is now a closed interval. If the filter had been [4000:4700], the entire second component would have been removed.

6 Header

The ASC Table Header contains metadata analogous to FITS header keywords. We allow ASC header attributes to have all the properties of a Descriptor, in contrast to FITS header keywords which do not have the full properties of a FITS table column.

6.1 Key Data Descriptor

A Key Data Descriptor has the same structure as a Table Column Data Descriptor. However, in the current implementation we will not support array dimensionality greater than 1 or axis group descriptors (cf. DSS Data Descriptor).

6.2 Grouping Descriptors

A new feature of the design is the idea of grouping descriptors. A group is a collection of related descriptors; its purpose is to allow software to display and select related information together.

By default, a block has only one group, the default group. It may have arbitrarily many named groups which collect together key, column, coordinate and filter descriptors. For instance, a 'time' group might include various timing keywords and the TIME table column. A group is NOT itself a kind of descriptor, and in particular group names do not need to be distinct from descriptor names.

Keys may be related to other 'parent' data descriptors, either other attributes or columns or data subspace axis descriptors. Attributes that are related to columns are called column attributes. Attributes that are related to data subspace axes are called data subspace attributes. All other attributes are table attributes. A generic FITS header keyword is a table attribute; the idea of tying header keywords to particular columns is new. A table attribute which is related to another table attribute may be considered as part of a group (equivalence class) of table attributes; this allows us to group header keywords and refer to them by groups rather than individually.

7 DM Images

7.1 Images and Tables

A DM Image is an DM Table with a single Table Column Data Descriptor whose array dimensionality $n > 0$ and with a single Row. Special access routines are provided for DM Images. Any single array Data Cell in a table may also be treated as a DM Image; to instantiate it as such an image, copy it to a new DM table together with the DSS, the Table Attributes, as well as the Data Descriptor and Column Attributes for its own Column, but discarding the other rows for the column and discarding the other column data descriptors, cells, and column attributes.

I illustrate the structure of an DM Image in the accompanying diagram; note that from the OO point of view this is just an instance of the DM Table, not a separate model.

8 Case studies and examples

8.1 FITS case study: PSPC off axis histogram file

An ASCII dump of a Rosat PSPC FITS file for the off axis histogram for an extracted source is reproduced below; I then interpret it in terms of the data model.

```
XTENSION= 'BINTABLE'           / binary table extension
BITPIX  =          8           / 8-bit bytes
NAXIS   =          2           / 2-dimensional binary table
NAXIS1  =          8           / width of table in bytes
NAXIS2  =         14           / number of rows in table
PCOUNT  =          0           / size of special data area
GCOUNT  =          1           / one data group (required keyword)
TFIELDS =          2           / number of fields in each row
TTYPE1  = 'OFF_AX_RAD'         / Off-axis grid point for histogram bin (arcmin)
TFORM1  = '1E'                / data format of the field: 4-byte REAL
TUNIT1  = 'arcmin'            / physical unit of field
TTYPE2  = 'FRAC_TIME'         / Fraction of time spent by source in bin
TFORM2  = '1E'                / data format of the field: 4-byte REAL
TUNIT2  = 'NONE'              / physical unit of field
EXTNAME = '0AH005'           / Detect extension-asp histogram for given source
```

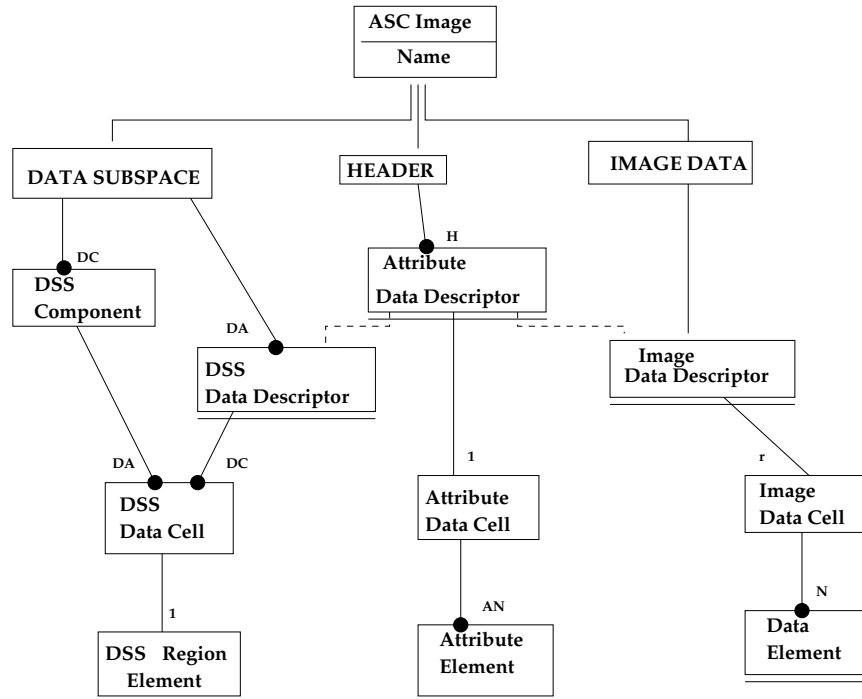


Figure 10: Data Model 7: DM Image Model, identical to Table Model but without Table Row and with only one Column Descriptor (Image Descriptor).

```

CONTENT = 'SOURCE ' / data content of file
ORIGIN = 'USRSDC ' / origin of processed data
DATE = '13/07/94' / FITS creation date (DD/MM/YY)
TELESCOP= 'ROSAT ' / mission name
INSTRUME= 'PSPCC ' / instrument name
OBS_MODE= 'POINTING' / obs mode: POINTING,SLEW, OR SCAN
IRAFNAME= 'rpi110590n00_oah005.tab' / IRAF file name
MJDREFI = 48043 / MJD integer SC clock start
MJDREFF = 8.79745370370074E-01 / MJD fraction SC clock start
ZERODATE= '01/06/90' / UT date of SC start (DD/MM/YY)
ZEROTIME= '21:06:50' / UT time of SC start (HH:MM:SS)
RDF_VERS= '2.9 ' / Rationalized Data Format release version number
RDF_DATE= '13-JUL-1994' / Rationalized Data Format release date
PROC_SYS= 'SASS7_2_0' / Processing system
PROCDATE= '2-JUN-1994 11:20:34' / SASS SEQ processing start date
REVISION= 2 / Revision number of processed data
FILTER = 'NONE ' / filter id: NONE OR BORON
OBJECT = 'XRT/PSPC PSF AR LAC' / name of object
RA_NOM = 3.320239E+02 / nominal RA (deg)
DEC_NOM = 4.551389E+01 / nominal DEC (deg)
ROLL_NOM= -1.349511E+02 / nominal ROLL (deg CCW North)
EQUINOX = 2.000000E+03 / equinox
OBS_ID = 'CA110590P.N10' / observation ID
  
```

```

ROR_NUM =          110590 / ROR number
OBSERVER= 'MPE, ROSAT-TEAM' / PI name
SETUPID = 'NOMINAL ' / Instrument setup
DATE-OBS= '20/06/90' / UT date of obs start (DD/MM/YY)
TIME-OBS= '11:24:43.000' / UT time of obs start (HH:MM:SS)
DATE_END= '20/06/90' / UT date of obs end (DD/MM/YY)
TIME_END= '13:07:12.000' / UT time of obs end (HH:MM:SS)
MJD-OBS =          4.806248E+04 / MJD of seq start
SCSEQBEG=          1606667 / SC seq start(sec)
SCSEQEND=          1612816 / SC seq end (sec)
NUM_OBIS=          2 / Number of obs intervals (OBIs)
LIVETIME=          1.884999E+03 / Live time
DTCOR =          9.602644E-01 / Dead time correction factor
ONTIME =          1.963000E+03 / On time
MPLSX_ID=          5 / Source number from merged source list (MPLSX)
EFFAREA =          1.0000E+00 / Effective area scaling factor
QUALITY =          0 / Quality of data (0 = good data)
RADECSYS= 'FK5 ' / WCS for this file
OFFAX =          1.478056E+01 / Off-axis angle of source in arcmin
COMMENT
COMMENT The following keywords are required in order to conform
COMMENT to the Office of Guest Investigator Programs standard:
COMMENT
AREASCAL=          1.0000E+00 / Area scaling factor
BACKFILE= 'NONE ' / No background file
BACKSCAL=          1.0000E+00 / Background scaling factor
CORRFILE= 'NONE ' / No correction file
CORRSCAL=          1.0000E+00 / Correction file scaling factor
RESPFILE= 'NONE ' / BLDRSP response file name (default)
ANCRFILE= 'NONE ' / BLDRSP ancillary response file name (default)
XFLT0001= ' ' / Required keyword
PHAVERSN= '1992A ' / Version # of OGIP file specification
POISSERR=          T / Poissonian errors appropriate
SYS_ERR =          0.0 / No systematic error
CHANTYPE= 'PI ' / Gain-corrected channels used
DETHANS=          256 / Total number of PHA channels available
COMMENT
COMMENT End required OGIP keywords
COMMENT
COMMENT This extension contains the off-axis histogram
COMMENT for the source given in the header.
HISTORY
HISTORY SASS file used: SPCBF.SEQ
HISTORY
HISTORY Correspondence with SASS variables:
HISTORY
HISTORY OFF_AX_RAD = OFF_SPB
HISTORY FRAC_TIME = OHS_SBP

OFF_AX_RAD  FRAC_TIME
0           0
5           0

```


10	2.17989E-04
15	0.73768
20	0.2621
25	0
30	0
35	0
40	0
45	0
50	0
55	0
57.5	0
60	0

What does this file contain? There's a lot of stuff all mixed together. We might describe it as follows:

Table OAH005(2 cols, 14 rows)

Colname	OFF_AX_RAD	FRAC_TIME
Datatype	Real(4)	Real(4)
Unit	none	none
Elt type	V	V
Elt dim	1	1
Disp	none	none
Desc	'Off-axis grid point for histogram bin (arcmin)' 'Fraction of time spent by source in bin'	
Component name	(same as colname)	
Array dim	0	0

Cells: 1 element per cell

Elements: 1 value per element (type V, dimension 1)

Values:

0	0
5	0
10	2.17989E-04
15	0.73768
20	0.2621
25	0
30	0
35	0
40	0
45	0
50	0
55	0
57.5	0
60	0

Data Subspace(4 axes)

TIME [1606667:1612816)

Coordinate: Origin = 0

Value = JD 2448044.379745370370074 d

```

Delta = 1
Unit = s
Comment SC seq start(sec)
Correction Factor 0.96026 (DTCOR)
RA/DEC Region not given (would be nice!)
2D Coordinate: Origin = not given
Value = J2000 (332.0239, +45.51389)

```

```

Delta = not given
Unit = deg

```

The following data subspace axes are not explicitly present in the file:

```

OFF_AX_RAD [0:60]
Unit = arcmin
Comment Off-axis grid point for histogram bin
FRAC_TIME [0:1]
Unit = none
Comment Fraction of time spent by source in bin

```

The following header cards from the file are not retained in our 'model' version as header cards per se because they contain information about the structure of the file or the attributes of its data axes:

\small

\begin{verbatim}

```

Cards from FITS standards, mapped to table structure:
XTENSION= 'BINTABLE' / binary table extension
BITPIX = 8 / 8-bit bytes
NAXIS = 2 / 2-dimensional binary table
NAXIS1 = 8 / width of table in bytes
NAXIS2 = 14 / number of rows in table
PCOUNT = 0 / size of special data area
GCOUNT = 1 / one data group (required keyword)
TFIELDS = 2 / number of fields in each row
TTYPE1 = 'OFF_AX_RAD' / Off-axis grid point for histogram bin (arcmin)
TFORM1 = '1E' / data format of the field: 4-byte REAL
TUNIT1 = 'arcmin' / physical unit of field
TTYPE2 = 'FRAC_TIME' / Fraction of time spent by source in bin
TFORM2 = '1E' / data format of the field: 4-byte REAL
TUNIT2 = 'NONE' / physical unit of field
EXTNAME = 'OAH005' / Detect extension-asp histogram for given source

```

Cards from OGIP rules, mapped to subspace and coordinate info:

```

MJDREFI = 48043 / MJD integer SC clock start
MJDREFF = 8.79745370370074E-01 / MJD fraction SC clock start
ZERODATE= '01/06/90' / UT date of SC start (DD/MM/YY)
ZEROTIME= '21:06:50' / UT time of SC start (HH:MM:SS)
RA_NOM = 3.320239E+02 / nominal RA (deg)

```

```

DEC_NOM =          4.551389E+01 / nominal DEC (deg)
ROLL_NOM=        -1.349511E+02 / nominal ROLL (deg CCW North)
EQUINOX =          2.000000E+03 / equinox
DATE-OBS= '20/06/90'          / UT date of obs start (DD/MM/YY)
TIME-OBS= '11:24:43.000'      / UT time of obs start (HH:MM:SS)
DATE_END= '20/06/90'          / UT date of obs end (DD/MM/YY)
TIME_END= '13:07:12.000'      / UT time of obs end (HH:MM:SS)
MJD-OBS =          4.806248E+04 / MJD of seq start
SCSEQBEG=          1606667 / SC seq start(sec)
SCSEQEND=          1612816 / SC seq end (sec)
LIVETIME=         1.884999E+03 / Live time
DTCOR  =          9.602644E-01 / Dead time correction factor
ONTIME  =          1.963000E+03 / On time

```

When writing this file back out, all of the above cards would be generated automatically by the FITS writing layer; there's no need for any of the software beyond the IO layer to ever deal with them.

The remaining header cards come in a number of groups, which we can't deduce from the present structure of the file:

Ungrouped header cards

```

OFFAX          14.78056
               Unit arcmin
               Comment Nominal off-axis angle of source

```

Header group PROCESSING

```

CONTENT = 'SOURCE'          / data content of file
ORIGIN  = 'USRSDC'          / origin of processed data
DATE    = '13/07/94'        / FITS creation date (DD/MM/YY)
IRAFNAME= 'rp110590n00_oah005.tab' / IRAF file name
RDF_VERS= '2.9'              / Rationalized Data Format release version number
RDF_DATE= '13-JUL-1994'     / Rationalized Data Format release date
PROC_SYS= 'SASS7_2_0'        / Processing system
PROCDATE= '2-JUN-1994 11:20:34' / SASS SEQ processing start date
REVISION=                    2 / Revision number of processed data

```

Header group OBSERVATION_DETAILS

```

TELESCOP= 'ROSAT'          / mission name
INSTRUME= 'PSPCC'          / instrument name
OBS_MODE= 'POINTING'        / obs mode: POINTING,SLEW, OR SCAN
FILTER   = 'NONE'           / filter id: NONE OR BORON
OBJECT   = 'XRT/PSPC PSF AR LAC' / name of object
OBS_ID   = 'CA110590P.N10'  / observation ID
OBSERVER= 'MPE, ROSAT-TEAM' / PI name
NUM_OBIS=                    2 / Number of obs intervals (OBIs)
ROLL_NOM=          -134.95

```

Header group ROSAT_SPECIFIC

```

ROR_NUM =          110590 / ROR number
SETUPID = 'NOMINAL'        / Instrument setup
MPLSX_ID=          5 / Source number from merged source list (MPLSX)
QUALITY  =          0 / Quality of data (0 = good data)

```

Header group OGIP_COMPAT / These keywords may be ignored by our software

```

EFFAREA =          1.0000E+00 / Effective area scaling factor
COMMENT
COMMENT  The following keywords are required in order to conform
COMMENT  to the Office of Guest Investigator Programs standard:
COMMENT
AREASCAL=          1.0000E+00 / Area scaling factor
BACKFILE= 'NONE    '          / No background file
BACKSCAL=          1.0000E+00 / Background scaling factor
CORRFILE= 'NONE    '          / No correction file
CORRSCAL=          1.0000E+00 / Correction file scaling factor
RESPFILE= 'NONE    '          / BLDRSP response file name (default)
ANCRFILE= 'NONE    '          / BLDRSP ancillary response file name (default)
XFLT0001= '        '          / Required keyword
PHAVERSN= '1992A  '          / Version # of OGIP file specification
POISSERR=          T          / Poissonian errors appropriate
SYS_ERR  =          0.0       / No systematic error
CHANTYPE= 'PI     '          / Gain-corrected channels used
DETCANS=          256       / Total number of PHA channels available
COMMENT
COMMENT  End required OGIP keywords
COMMENT

```

Header group COMMENTS

```

COMMENT  This extension contains the off-axis histogram
COMMENT  for the source given in the header.
HISTORY
HISTORY  SASS file used: SPCBF.SEQ
HISTORY
HISTORY  Correspondence with SASS variables:
HISTORY
HISTORY  OFF_AX_RAD = OFF_SPB
HISTORY  FRAC_TIME  = OHS_SBP

```

How would we redesign this file to take more advantage of the data model while remaining compatible with software that expects the old format? While I do not expect that we will be writing software to regenerate PSPC standard data products in this way, it's a useful exercise to show what is needed to add the extra structure.

- We add comments to denote Header Groups, grouping the table attributes. This could be used by browsers to organize the user's view of the data. It would be nice for software to be able to use such header groups, but there is a risk that some FITS readers will mangle the order of the header keywords, mixing up the group memberships. I still feel that it's an enhancement worth having, with the warning to users that if they pass the files through other software they may lose that information.
- The other way of making header groups is to explicitly add named cards. This is comparatively inefficient but may be the way to go when it's important that the linkage be robust. This is illustrated with the DAREL keywords for OFFAX and ROLLNOM.
- The dataset is actually binned data; the OFF_AX_RAD column contains bins which for some perverse reason are uneven in size near the ends. I could have defined a special element type to denote bins where the boundaries are deduced to be half way to the next entry, but this would require the software to handle more than one row at a time. I prefer to accept the overhead of the extra two columns COL1_LO and COL1_HI, turning OFF_AX_RAD into a column of element type T (two sided uncertainty).

- We will store the extraction region in the data subspace header. The information includes the region specification in sky pixel coordinates and the transformation from sky pixel coordinates to RA and Dec, the latter being copied from the original file. This gives us a more logical place to put the info now stored in RA_NOM and DEC_NOM. If we had the region specification in RA and Dec instead of pixels, we would store it in keyword DSC1 instead of DS1.
- The preferred columns are OFF AX RAD and TIME; but we don't need to include PREF1 and PREF2 keywords since these are the only two columns at the data model level and they are in the correct order.

```

XTENSION= 'BINTABLE'          / binary table extension
BITPIX   =                    8 / 8-bit bytes
NAXIS    =                    2 / 2-dimensional binary table
NAXIS1   =                   16 / width of table in bytes
NAXIS2   =                   14 / number of rows in table
PCOUNT   =                    0 / size of special data area
GCOUNT   =                    1 / one data group (required keyword)
TFIELDS  =                    4 / number of fields in each row
TTYPE1   = 'OFF_AX_RAD'      / Off Axis Radius
TFORM1   = '1E'              / data format of the field: 4-byte REAL
TUNIT1   = 'arcmin'         / physical unit of field
TTYPE2   = 'COL1_LO'        / Lower Uncertainty
TFORM2   = '1E'              / 4 byte real
TUNIT2   = 'arcmin'         /
TTYPE3   = 'COL1_HI'        / Upper Uncertainty
TFORM3   = '1E'              / 4 byte real
TUNIT3   = 'arcmin'         /
TTYPE4   = 'FRAC_TIME'      / Fractional Exposure Time
TFORM4   = '1E'              / data format of the field: 4-byte REAL
TUNIT4   = '                ' / physical unit of field
EXTNAME  = 'OAH005'         / Off Axis Histogram
TDISP1   = 'F8.2'           / Format to display OFF AX RAD
TDISP4   = 'F8.6'           / Format to display FRAC TIME
TLMIN1   =                   0.0 / Valid range for columns
TLMAX1   =                   60.0 /
TLMIN4   =                   0.0 /
TLMAX4   =                   1.0 /
COMMENT
COMMENT ASC Table Keywords
COMMENT
DCFIELDS=                    2 / Number of logical columns
DCETYP1  = 'T'               / Two sided uncertainty
DCITYP1  = '['               / Interval type
COMMENT
COMMENT ASC Data Subspace Keywords
COMMENT
DSNAXIS  =                    1 / Number of data subspace axes
DSNAM1   = 'SKYPOS'         / Sky pixel position
DSDIM1   =                    2 / Dimension of DSNAM1
DSTYP1   = 'X'              / First component of DSNAM1
DSTYP2   = 'Y'              / Second component of DSNAM1
DSUNIT1  = 'pixel'          /
DSCNAM1  = 'EQPOS'          / Coordinate system on DSNAM1

```

```

DSCTYP1 = 'RA---TAN'           / Transform for axis 1
DSCTYP2 = 'DEC--TAN'          / Transform for axis 2
DSCUNI1 = 'deg'               /
DSCRVL1 =          332.0239    / Reference RA value (RA\_NOM)
DSCRVL2 =          45.5138    / Reference Dec value (DEC\_NOM)
DSCRFX1 =          4096.5000   / Reference X value
DSCRFX2 =          4096.5000   / Reference Y value
DSCDLT1 =          -0.0124    / Deg per pixel
DSCDLR2 =           0.0124    / Deg per pixel
DS1      = 'c 4087.3 4012.3 43.2' / Extraction region in X,Y coords
DSTYP3   = 'TIME'             / Mission time
DSUNIT3  = 's'                /
DS2L1    =          1606667.0  / Start time
DS2U1    =          1612816.0  / Stop time
DSITYP3  = '[]'              / Interval type for TIME
COMMENT
COMMENT Alternative syntax for the above three keywords would be:
COMMENT   DS2 = '[SCSEQBEG:SCSEQEND]'
COMMENT
DSCTYP3  = 'DATE'            / Calendar date
DSCDLT3  =          1.15741E-05 / Days per second
DSCRVL3  = 48043.879745370370074 / MJD of SC clock start
DSCRFX3  =           0.0       / SC clock start
DSCUNI3  = 'd'               /
DSTYP4   = 'OFF_AX_RAD'      / Range defaults to TLMIN1/TLMAX1
DSTYP5   = 'FRAC_TIME'       /
COMMENT
COMMENT ASC Table Attributes
COMMENT
COMMENT We only need to use explicit DANAMn keywords when we
COMMENT want to add extra information to a keyword.
COMMENT

DANAM1   = 'OFFAX'           / Attribute
OFFAX    =          1.478056E+01 / Off-axis angle of source in arcmin
DAUNI1   = 'arcmin'          / Unit of DANAM1
DAREL1   = 'OFF_AX_RAD'      / Keyword OFFAX is bound to column OFF AX RAD

DANAM2   = 'ROLL_NOM'        / Attribute
ROLL_NOM=          -1.349511E+02 / nominal ROLL (deg CCW North)
DAUNI2   = 'deg'            /
DAREL2   = 'SKYPOS'          / ROLL_NOM bound to DSS axis SKYPOS

DANAM3   = 'SRC_OFF_AX_RAD'   / Same as OFFAX,
DAVAL3   =          1.478056E+01 / but illustrating a name longer than 8 chars
DAUNI3   = 'arcmin'          /

DANAM4   = 'ONTIME'          / Denote the fact that the keywords named
DAREL4   = 'TIME'            / are tied to the TIME information, so if that
DANAM5   = 'DTCOR'           / becomes invalid so do these.
DAREL5   = 'TIME'            / Debatable whether we would actually bother
DANAM6   = 'LIVETIME'        / to add these linkages in this case.

```

```

DAREL6 = 'TIME ' /

COMMENT
COMMENT Header Group PROCESSING
COMMENT
CONTENT = 'SOURCE ' / data content of file
ORIGIN = 'USRSDC ' / origin of processed data
DATE = '13/07/94' / FITS creation date (DD/MM/YY)
IRAFNAME= 'rp110590n00_oah005.tab' / IRAF file name
RDF_VERS= '2.9 ' / Rationalized Data Format release version number
RDF_DATE= '13-JUL-1994' / Rationalized Data Format release date
PROC_SYS= 'SASS7_2_0' / Processing system
PROCDATE= '2-JUN-1994 11:20:34' / SASS SEQ processing start date
REVISION= 2 / Revision number of processed data
COMMENT
COMMENT Header Group Observation Details
COMMENT
TELESCOP= 'ROSAT ' / mission name
INSTRUME= 'PSPCC ' / instrument name
OBS_MODE= 'POINTING' / obs mode: POINTING,SLEW, OR SCAN
FILTER = 'NONE ' / filter id: NONE OR BORON
OBJECT = 'XRT/PSPC PSF AR LAC' / name of object
OBS_ID = 'CA110590P.N10' / observation ID
OBSERVER= 'MPE, ROSAT-TEAM' / PI name

COMMENT
COMMENT Header Group ROSAT Specific
COMMENT

ROR_NUM = 110590 / ROR number
SETUP ID = 'NOMINAL ' / Instrument setup

COMMENT
COMMENT Header Group HEASARC Position Keywords
COMMENT

RA_NOM = 3.320239E+02 / nominal RA (deg)
DEC_NOM = 4.551389E+01 / nominal DEC (deg)
EQUINOX = 2.000000E+03 / equinox
RADECSYS= 'FK5 ' / WCS for this file

COMMENT
COMMENT Header Group HEASARC Timing Keywords
COMMENT
MJDREFI = 48043 / MJD integer SC clock start
MJDREFF = 8.79745370370074E-01 / MJD fraction SC clock start
ZERODATE= '01/06/90' / UT date of SC start (DD/MM/YY)
ZEROTIME= '21:06:50' / UT time of SC start (HH:MM:SS)
DATE-OBS= '20/06/90' / UT date of obs start (DD/MM/YY)
TIME-OBS= '11:24:43.000' / UT time of obs start (HH:MM:SS)

```

```

DATE_END= '20/06/90'           / UT date of obs end   (DD/MM/YY)
TIME_END= '13:07:12.000'      / UT time of obs end  (HH:MM:SS)
SCSEQBEG=          1606667 / SC seq start(sec)
SCSEQEND=          1612816 / SC seq end  (sec)
MJD-OBS =          4.806248E+04 / MJD of seq start
NUM_OBIS=           2 / Number of obs intervals (OBIs)
LIVETIME=          1.884999E+03 / Live time
DTCOR  =           9.602644E-01 / Dead time correction factor
ONTIME  =           1.963000E+03 / On time
MPLSX_ID=           5 / Source number from merged source list (MPLSX)
EFFAREA =           1.0000E+00 / Effective area scaling factor
COMMENT
COMMENT Header Group OGIP_COMPAT
COMMENT
COMMENT The following keywords are required in order to conform
COMMENT to the Office of Guest Investigator Programs standard:
COMMENT
AREASCAL=           1.0000E+00 / Area scaling factor
BACKFILE= 'NONE'      ' / No background file
BACKSCAL=           1.0000E+00 / Background scaling factor
CORRFILE= 'NONE'      ' / No correction file
CORRSCAL=           1.0000E+00 / Correction file scaling factor
RESPFILE= 'NONE'      ' / BLDRSP response file name (default)
ANCRFILE= 'NONE'      ' / BLDRSP ancillary response file name (default)
XFLT0001= ' '          ' / Required keyword
PHAVERSN= '1992A'     ' / Version # of OGIP file specification
COMMENT Note that the error info given here applies to the counts errors
COMMENT which are in an entirely different table; so we don't
COMMENT attach them to the data model errors in this file.
POISSERR=           T / Poissonian errors appropriate
SYS_ERR =           0.0 / No systematic error
CHANTYPE= 'PI'       ' / Gain-corrected channels used
DETHANS=           256 / Total number of PHA channels available
COMMENT
COMMENT End required OGIP keywords
COMMENT
COMMENT Header Ungrouped
COMMENT

QUALITY =           0 / Quality of data (0 = good data)
COMMENT
COMMENT This extension contains the off-axis histogram
COMMENT for the source given in the header.
HISTORY
HISTORY SASS file used: SPCBF.SEQ
HISTORY
HISTORY Correspondence with SASS variables:
HISTORY
HISTORY OFF_AX_RAD = OFF_SPB
HISTORY FRAC_TIME  = OHS_SBP

```


OFF_AX_RAD	COL1_LO	COL1_HI	FRAC_TIME
0	0	2.5	0
5	2.5	2.5	0
10	2.5	2.5	2.17989E-04
15	2.5	2.5	0.73768
20	2.5	2.5	0.2621
25	2.5	2.5	0
30	2.5	2.5	0
35	2.5	2.5	0
40	2.5	2.5	0
45	2.5	2.5	0
50	2.5	2.5	0
55	2.5	1.25	0
57.5	1.25	1.25	0
60	1.25	0	0

8.2 Case Study: Barycenter Correction Algorithm

We analysed the Barycenter Correction Algorithm to see how it would be laid out in terms of the data model.

The algorithm uses the following ASC Tables:

- **Event List:** this contains rows which we refer to as photons, and a set of columns which include at least **Pixel Position** and **Time**. The Pixel Position Column Data Descriptor has Data Descriptor with default name Pixel Position and component names X and Y; it must be of element dimension 2. We will access it by element type V (Value). It must also have a Data Coordinate Descriptor which contains the Equatorial Position (RA and Dec). The Time Data Descriptor may have a Data Coordinate Descriptor giving the absolute Date.
- **Orbital Data:** This is a stack containing the names of spacecraft and pointers to their Ephemeris files.
- **Solar System Ephemeris:** This is a stack containing the names of planets and pointers to their Ephemeris files.
- **Ephemeris:** This is a table with the columns Time and 3-Vector-Position. The latter has element dimension 3 and component names X,Y,Z. The ephemeris table has a table attribute Mass, giving the mass of the orbiting body.

The algorithm is:

- Identify the spacecraft in use for this event list: this should be a table attribute of the event list.
- Find the corresponding spacecraft ephemeris from the orbital data stack.
- Open output table with same format as input event list but with extra column named BARY_TIME of dimension 1 and type U. Unit is seconds of mission time; coordinate system is copied from input column whose default name is TIME. Add comment to header describing the fact that BARY_TIME is the time of a different event (arrival of a photon at the barycenter) in the same coordinate system as TIME.
- For each row in the table, get the Pixel Position. Calculate the Equatorial Position using that Data Descriptor's Data Coordinate Transform.
- Calculate the 3-vector direction of the photon (the source vector) from the equatorial position.
- Get time from row of table (represents photon arrival time at spacecraft). (If the ephemerides are in JD rather than mission time, may need to also use this Data Descriptor's Data Coordinate Transform to get JD from time.) Also get time uncertainty if present.
- Interpolate in spacecraft ephemeris at the given time to return the spacecraft ephemeris position and uncertainty (an element of type U and dimension 3).
- For each entry (planet) in the solar system ephemeris stack, interpolate in the corresponding ephemeris and return the mass of the planet and the position (a value element of dimension 3) at the time.
- Calculate the solar system barycenter at the given time by taking the mass weighted mean of the planetary positions. Result is an element of type V and dimension 3.

- Calculate the barycenter to spacecraft vector and its uncertainty. Check that the units of barycenter and spacecraft positions are compatible and apply conversions if necessary.
- Calculate the scalar product of the spacecraft and source vectors and its uncertainty; scale to light travel time to obtain correction. Correction is an element of type U.
- Add this to photon time and combine uncertainty in quadrature. Result is barycenter corrected time (BARY_TIME).
- Copy input row to output, adding new column of BARY_TIME.
- Loop to next photon until complete.

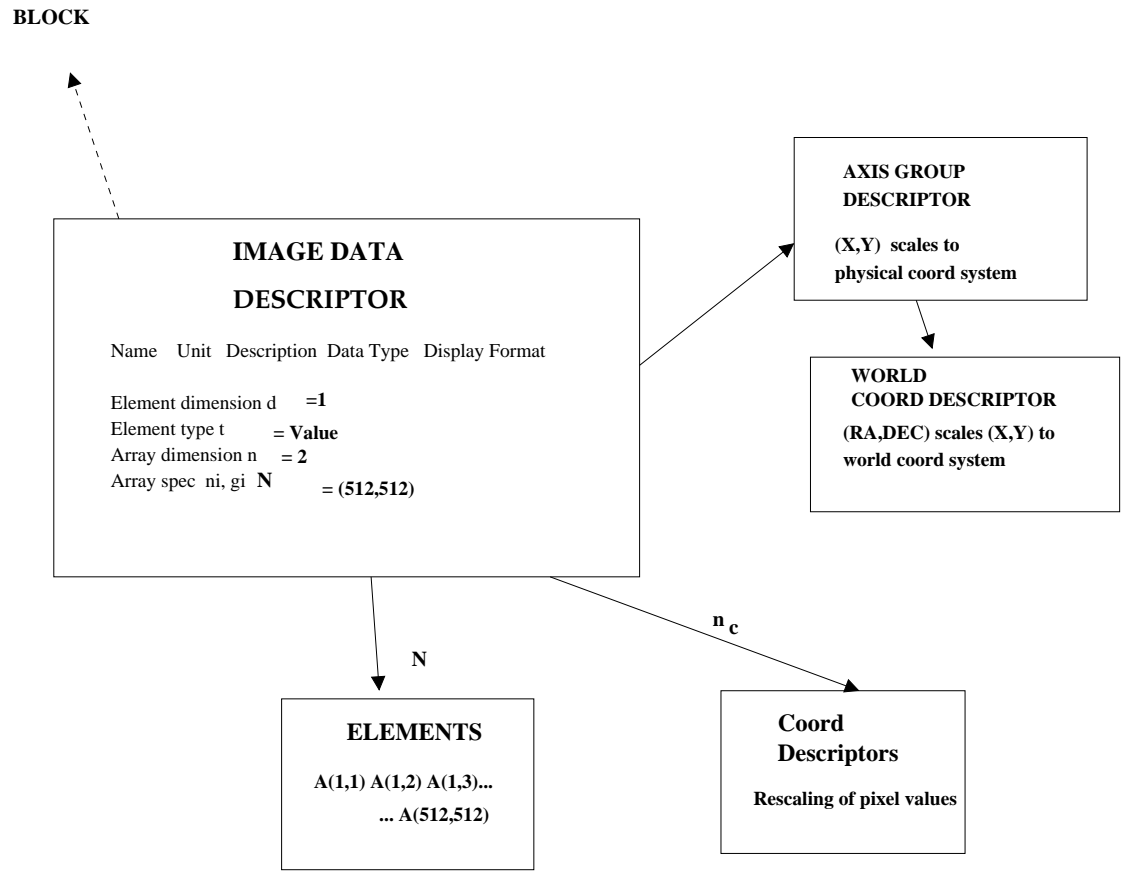


Figure 11: Data Model 7b: DM Image Descriptor