



International

Virtual

Observatory

Alliance

Simulation Data Model

Version 1.00-20111019

IVOA DM WG and TIG Proposed Recommendation 2011 October 19

This version:

1.00-20111019

Latest version:

1.00-20111019,

<http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/PR-SimulationDataModel-v.1.00-20111019.doc>

Previous version(s):

See revision page on GoogleCode:

<http://code.google.com/p/volute/source/browse/trunk/projects/theory/snapdm/specification/>

Working Group:

<http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/IvoaDataModel>

<http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/IvoaTheory>

Editors:

Gerard Lemson, Hervé Wozniak

Authors:

Gerard Lemson, Laurent Bourgès, Miguel Cerviño, Claudio Gheller, Norman Gray, Franck LePetit, Mireille Louys, Benjamin Ooghe, Rick Wagner, Hervé Wozniak

Abstract

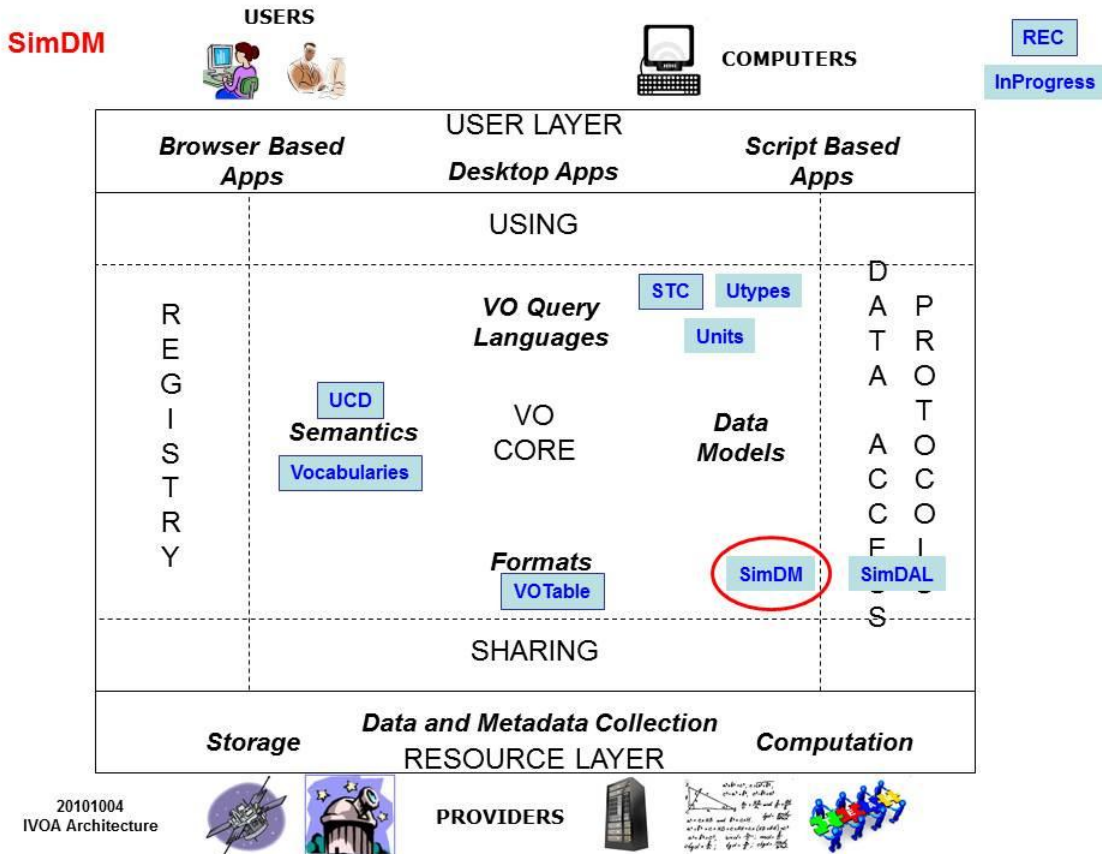
In this document and the accompanying documents we propose a data model (Simulation Data Model) describing numerical computer simulations of astrophysical systems. The primary goal of our proposal is to support discovery of simulations by describing those aspects of them that scientists might wish to query on, i.e. it is a model for *meta*-data describing simulations.

This document does *not* propose a protocol for using this model. Two distinct IVOA protocols (SimDB and SimDAL) are in the make and both are supposed to use the model, either in its original form or in a form derived from the model proposed here, but more suited to the particular protocol.

The SimDM has been developed in the IVOA Theory Interest Group with assistance of representatives of relevant working groups, in particular DM and Semantics.

Link to IVOA Architecture

The figure below shows where SimDM fits within the IVOA architecture:



Status of This Document

This is an IVOA Proposed Recommendation made available for public review. It is appropriate to reference this document only as a recommended standard that is under review and which may be changed before it is accepted as a full recommendation.

The first release of this document was 2011 April 28.

A list of current IVOA Recommendations and other technical documents can be found at <http://www.ivoa.net/Documents/> .

Acknowledgements

We thank various persons for useful discussions in the course of this work: first, the participants of the [Feb 2006 theory workshop](#) in Cambridge, UK, where this work was started; second, the participants of the [April 2007 SNAP workshop](#) in Garching, Germany, where the design started taking shape. The work has also been influenced by the participants of the Technical Coordination Group of the EuroVO-DCA project and participants of the theory workshop organised in the context of that project in [Garching, 2008](#). Then we want to thank particularly the following persons for useful discussions and feedback: Jeremy Blaizot, Miguel Cerviño, Klaus Dolag, Pierro Madau, Adi Nusser, Ray Plante, Volker Springel, and Alex Szalay. We finally want to thank participants to the theory sessions in all the interoperability meetings since Victoria 2006, where parts of this work were discussed.

Conformance related definitions

The words "MUST", "SHALL", "SHOULD", "MAY", "RECOMMENDED", and "OPTIONAL" (in upper or lower case) used in this document are to be interpreted as described in IETF standard, RFC 2119 [0].

The **Virtual Observatory (VO)** is a general term for a collection of federated resources that can be used to conduct astronomical research, education, and outreach. The **International Virtual Observatory Alliance (IVOA)** is a global collaboration of separately funded projects to develop standards and infrastructure that enable VO applications. The International Virtual Observatory (IVO) application is an application that takes advantage of IVOA standards and infrastructure to provide some VO service.

Contents

1	Introduction	5
2	SimDM: application, approach and outline	7
2.1	Application(s) of the model	7
2.2	Modelling approach	8
2.3	Phase 1: analysis	9
2.4	Phase 2: domain model	10
3	Logical model overview	13
3.1	Packages	14
3.2	Resource	14
3.3	Physics, models and algorithms	17
3.4	Parameters: definition and values	19
3.5	Target: Goal of experiment	20
3.6	Object types: real and simulated	22
3.7	Results: data sets and their statistical summary	23
3.8	Data access services	27
4	Serialisations	28
4.1	SimDM/UTYPE	28
4.2	XML	30
5	Dependencies on other IVOA efforts	31
5.1	Registry	31
5.2	Semantics: Use of SKOS Concepts	33
5.3	Data Model	34
5.3.1	UML Profile	34
5.3.2	Characterisation data model	34
5.3.3	UTYPE	35
6	References	35
6.1	Accompanying documents	35
6.2	Relevant IVOA documents	36
6.3	Other sources	37

1 Introduction

In this document we make a proposal for an IVOA standard data model for describing simulations¹. Indeed, apart from limited support for publishing model spectra in SSAP, there is as yet no IVOA standard dealing with the publication of simulations and their results. The primary goal of our proposal is to support discovery of simulations by describing those aspects of them that scientists might wish to query on, i.e. it is a model for *meta*-data describing simulations. This document does *not* propose a protocol for using this model. Two distinct IVOA protocols are in the make and both are supposed to use the model, either in its original form or in a form derived from the model proposed here, but more suited to the particular protocol.

The direct motivation of the model comes from the Simulation Database (SimDB) and the former Simulation Data Access Protocol (SimDAP) efforts. Work on these standards started under the header *Simple Numerical Access Protocol* (SNAP) in the Theory Interest Group (TIG) since the Victoria interoperability meeting, 2006. It was agreed in the Trieste interop 2008 to split SNAP in two separate tracks, SimDB and SimDAP. More recently it was deemed useful to further split the work on SimDB in two tracks, one focusing on the data model alone, the second on its serialisation and usage in the SimDB protocol. This *Note* deals with the data model, referred to as SimDM.

Work on the SimDB specification has been organised via a GoogleCode SVN repository in the *volute* project originally created by Norman Gray for the Semantics Working group. The history of the SimDB project can be obtained from <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm>. A large part of that though deals with technical issues revolving around a code generator we designed to derive relevant resources from the basic data model. Most of that development has been moved to a separate GoogleCode project, VO-URP. The resources dealing with the SimDB developments have been gathered in the subdirectory <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/> in <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/>. The other parts of the volute part of SimDB should be deemed deprecated.

The other documents related to this proposal and being a part of the specification are the following. They can be found at the root URL <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/> **[TO BE MOVED TO IVOA REPOSITORY]**:

¹ We will use the term *simulations* for the running of a simulation code as well as for their results. And we will often include post-processing codes and their results as well.

Table 1: list of documents as part of the SimDM specification, accessible from <http://ivoa.net/XXX> [TO BE COMPLETED AFTER APPROVAL]

SimDM.html	Full browsable specification of the model	html/SimDM.html
SimDM_DM.png	Graphic view of the whole model (large image)	uml/SimDM_DM.png
SimDM_DM.xml	MagicDraw UML diagram serialised to XMI	uml/SimDM_DM.xml
SimDM_INTERMEDIATE.xml	Intermediate representation of the model: a (generated) XML document representing the complete model in more readable format than XMI	uml/SimDM_INTERMEDIATE.xml
intermediateModel.xsd	XML schema document for intermediate representation's XML format	uml/intermediateModel.xsd
xsd/	XML schema documents (generated) representing mapping of UML to XSD	xsd/

Additionally, an implementation note explains how to use this model to describe various kinds of theoretical products.

Section 2 described our methodology whereas the model itself is detailed in Section 3. A few specific issues of serialization are addressed in Section 4. The development of SimDM is linked to other IVOA efforts that deserve to be mentioned. Section 5 deals on that point. In the Appendices we deal with the scientific motivation at the origin of creating SimDM and the 4-years history of the developments (Appendix A), more details on the UML profile (Appendix B), some issues (Appendix C) and the use of the model in two cases (SimDB and SimDAL, see Appendix D).

2 SimDM: application, approach and outline

2.1 Application(s) of the model

The data model proposed here was never meant to be created in a vacuum (see Appendix A for more details), but was always intended to be used in some IVOA standard service. What precisely this application was has not always been clear for the SimDM, the goal of which has gone through some changes in the course of the project. It started off as a SIAP-like service protocol for N-body simulations, SNAP². Then mesh simulations were included, leading to our definition that SNAP should support the discovery and possibly partial retrieval of simulations involving “evolving objects in 3D space”³. Over time the protocol separated into the Simulation Database, a more TAP service for querying for simulations, and the SimDAP protocol for retrieving actual data products⁴. Recently (and finally) also simulations of a different type, “micro-physic simulations”⁵, or “models”⁶ have been included, and SimDAP has merged with S3 to form SimDAL: a family of access protocols for theory data⁷.

Consequently the aim and possible applications of the data model have changed a little over time. The discussions in the Victoria 2010 interop have finally led to a convergence of these ideas and to the following agreement:

1. The SimDM MUST support the SimDB protocol and the SimDAL service protocols.
2. To do so it MUST allow scientists to describe their simulations in sufficient detail for others to decide whether a given simulation is of interest, and to query for these simulations.

² <http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/CambridgeTheoryWorkshopFeb06>

³ <http://www.ivoa.net/internal/IVOA/InterOpMay2006Theory/closingplenary.ppt>

⁴ <http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/InterOpMay2008Theory>

⁵ Coined in Cambridge interop, 2007.

⁶ In Victoria 2010, we decided to label all as simulations. See http://www.ivoa.net/internal/IVOA/InterOpMay2010Theory/IVOA2010_ModelvsSim.pdf

⁷ <http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/InterOpMay2010Theory>

3. It SHOULD then offer directions to services for further drilling down or downloading simulations or parts of these.
4. The precise representation of the model in the individual service protocols (SimDB, SimDAL) is not prescribed.
5. If individual protocols will deviate in the details from the SimDM, this will be for application specific reasons.
6. The SimDM will provide the vocabulary for the concepts used in these possible alternative representations.

In the preferred case, these other representations can be considered as different *views* on the core SimDM. This should be interpreted similar to the way relational database views provide a mechanism by which slightly different representations from the often more normalised data model can be provided.

2.2 Modelling approach

With this application area in mind we have followed a somewhat formal approach to our data modelling effort, suggested by e.g. [27]. In this approach the construction of a data model is divided in three stages. The first is the *analysis* stage, in which one investigates the “domain of the application” one tries to model. This stage produces a *conceptual*, or *domain model* [32] which is relatively abstract and high level. In its design one does *not* aim to create a model directly suited to an application. The emphasis is on identifying important concepts and their relationship in the “real world”. One also refrains from giving a fully detailed definition of all attributes and other features. Important in this stage is interaction with *domain experts*, in our case scientists.

The aim of the second stage is the creation of a *logical model* [33] of the application domain. This should be a detailed model supporting the application. It should contain all information required to support the application. It should however still be implementation neutral, and concentrate on describing the precise concepts and semantic relationships between them. In general it will use part of the concepts from the domain model, but works them out in more, all, detail.

In the third stage one derives from this logical model one or more *physical models* [34]. These are representations of the logical model in a form that can be used directly by the various computational components building the system. Examples of this are XML schemas defining valid XML documents, or relational schemas for the design of a database storing instances of the model.

The analysis phase is described in the next few subsections. The logical model is defined in full detail in UML using the MagicDraw Community Edition modelling tool. MagicDraw stores the model in an XML file following the XMI serialisation [29] we have created a fully cross-linked HTML file documenting this model in all detail. These 2 files are part of this specification and can (for now⁸) be found on the GoogleCode volute site.

⁸ Once the specification becomes a proposed recommendation the stable versions of these files should be placed under the IVOA wiki site.

2.3 Phase 1: analysis

The analysis phase investigates the “world the application lives in”, its “universe of discourse” [27] and describes it in a domain model. To get constraints on this universe and its contents we follow [35] in trying to gather some 20 science questions that the application should be able to answer. The application is here a system consisting of the data model together with the protocol and implementations. The model will be designed in such a way that it can contain the required information. The protocol and implementations must support efficient querying for this information. [35] used this approach in the design of the SDSS database.

To create such a list of questions we have contacted scientists with the question that if they were presented with a database of simulation metadata, what questions they would want to ask of it to find interesting simulations. The following list summarises their answer:

- What system/object is being simulated?
- What physical processes are included?
- How is the system being represented in the simulation (particles (Lagrangian), (adaptive) mesh (Eulerian)), both, other?
- How are the physical processes implemented?
- What numerical approximations were used (e.g. resolution, softening parameter)?
- What observables are available for the system/object, possibly as function of time⁹? As it is a spatial system, at least *simulation boxsize*, centre-of-mass position.
- What observables are available for the constituents, i.e. what is the schema of the objects from which the simulation built e.g. particles in N-body simulation, grid cells in an adaptive mesh simulation or particle groups in a cluster finder?
- Per snapshot, per simulation object type, per variable:
 - Characterise the *possible* values
 - Characterise the result
- Are post-processing results available?
- Are services/applications available for accessing the results?
- Which code ran the simulation?
 - Which *version* of the code?
 - Is software available?
- Who ran the simulations?
- What were values of input parameters?
- How were initial conditions created?
- How the results are parameterized?
- Can I access grids of models? Can I access individual results?

⁹ Re: Rick Wagner’s example of certain properties only being calculated after a certain stage in the simulation is reached.

- Which are the inputs ingredients (usually, which data collections are used?)
- How I can run a simulation? Can I do it on-the-fly?
- Can include my simulations in the VO in an easy way? What I should do?
- Can i compare different simulations? Can I compare the simulation with my data?
- Which simulations provide diagnostic tools? (i.e. distance/extinction/quasi-scale free quantities)
- Can I combine the results of different simulations in a single file adapted for my needs (e.g. own code)?

2.4 Phase 2: domain model

The result of the analysis phase is a model in its own right, albeit rather sparse and schematic. For this purpose we have built on previous work by adapting the so called *Domain model for Astronomy* proposed in [11]. This model forms the basic structure of the domain model for SimDB, illustrated in Figure 1.

Figure 1 is used in a narrative motivating the final structure of the full SimDM. We start by assuming the existence of one or more **Files** that a publisher thinks may be of interest to the community because they contain astronomical data. Instead of in files the data might also reside in a **Database**, and to be generic we introduce a **Storage** base class that abstracts the actual physical location of the data.

Registering that files exist somewhere is not of great interest without providing information about the *contents* of the files. The philosophy that we follow is that the files are of potential interest because they contain the **Results**¹⁰ of an (astronomical) **Experiment**, and accordingly their contents must be explained by describing the experiment that gave rise to it. Only in this way can one make scientific use of the files or other storage resources.

The abstract **Experiment** is made concrete by adding some examples of experiment types that are important for the current model dealing with **Simulations** and simulation **PostProcessing**.

In our model, **Experiment** represents the actual *running* of an experiment; to describe the *design* of the experiment (the so-called *experimental protocol*) we introduce the concept of (*experimental*) **Protocol**¹¹. This separation between design of experiment and the execution is a *normalisation* that reduces redundancy in the model. See the accompanying appendix for a discussion of this concept. We mirror the concrete subclasses of **Experiment** by adding

¹⁰ We do not assume that in reality the relation between the conceptual Result and the concrete Storage elements can be modelled by a single reference. Especially for the largely non-standardised world of simulations a single result can be distributed over many files, but it is also possible for one file to contain multiple results. In the current SimDB model we do not attempt to model such relations explicitly. We delegate the responsibility for accessing the physical results to (web) services and this issue is more explicitly addressed by the SimDAP protocol.

¹¹ Further on, the word *protocol* will be preceded by the adjective *experimental* (in parenthesis and italicized) to keep clear the distinction with any other IVOA protocol.

concrete subclasses to *(experimental)* **Protocol** such as **Simulator**, which represents simulation codes according to which **Simulations** are run, and **PostProcessor** corresponding to **PostProcessing** runs.

The *(experimental)* **Protocol** class contains **InputParameters**. An **Experiment** using a particular *(experimental)* **Protocol** only needs to indicate the *values* for these parameters. In this way a single instance of the *(experimental)* **Protocol** can be reused by many **Experiments** performed according to it.

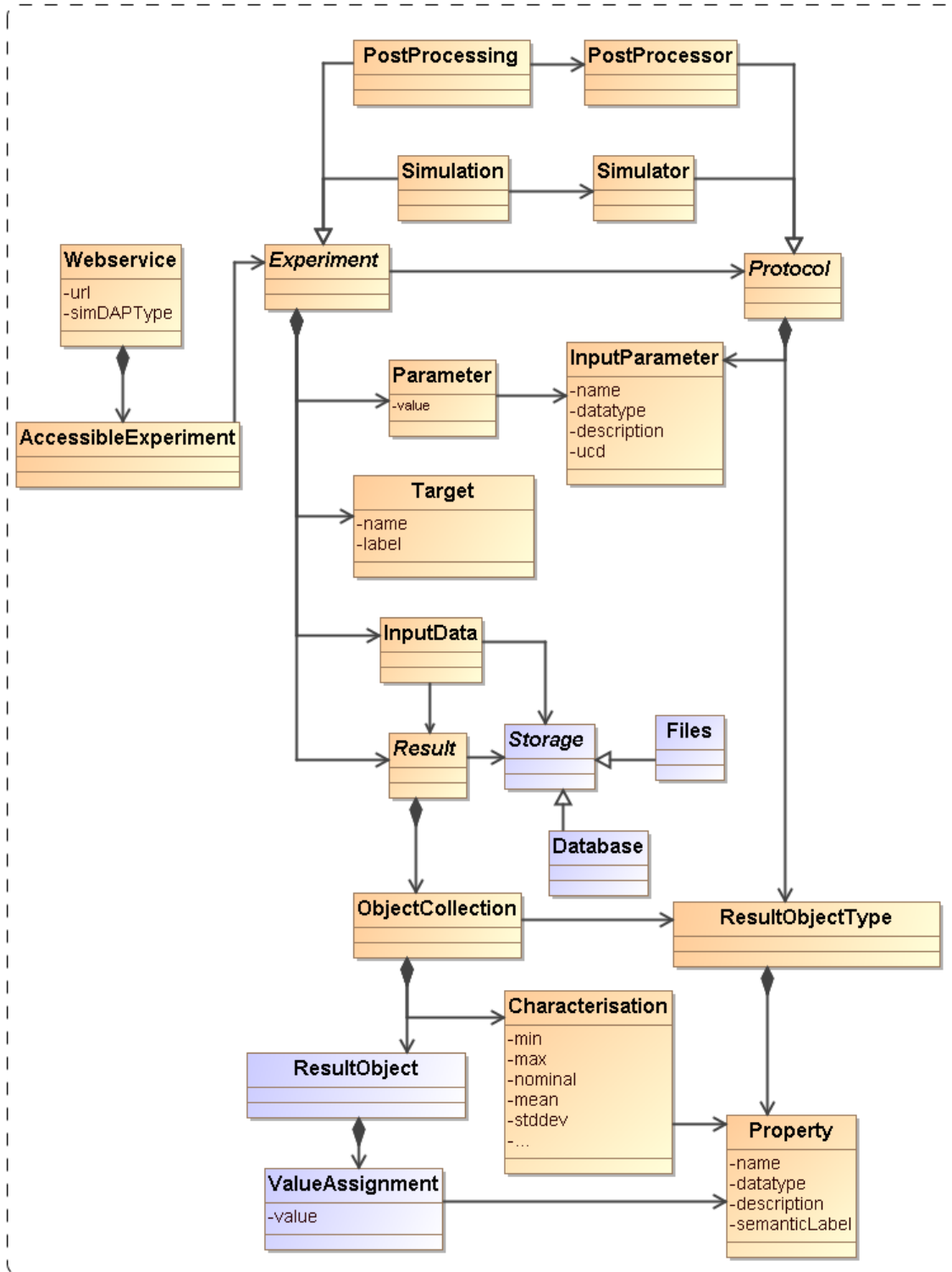


Figure 1: Schematic domain model encapsulating the main design constructs in SimDM. Elements coloured orange are represented directly in SimDM, possibly with a different name. Purple elements are not part of that model, but are used to explain and motivate other features that do appear there.

The (*experimental*) **Protocol** also defines the possible structure of the results of the experiments. In our model **Results** contain **ResultObjects**. These objects have a given type, represented by the **ResultObjectType** contained by (*experimental*) **Protocol**. The **ResultObjectType** defines the **Properties** that these objects have.

For example the results of N-body simulations may contain particles having properties position, velocity, mass and possibly others. Adaptive Mesh Refinement (AMR) simulations produce results that are collections of mesh cells of various sizes, positions and contents. Similarly post-processing codes such as halo finders produce “halos” and “semi-analytical” galaxy formation codes produce galaxies.

In general a single result can contain objects of different types. For example a Smooth Particle Hydrodynamics (SPH) simulation may contain dark matter particles, star particles and gas particles. And in general the codes allow one to configure which of these exactly are chosen in a given experiment.

One aspect of the experiment that is not determined by the experimental protocol is *why* the experiment was performed. In the model we introduce the **Target** concept for this, which represents real world objects or processes that are being simulated. For example, with the same N-body simulator one may simulate a galaxy merger or the evolution of large scale structure of the universe.

As discussed above, the actual way in which results are stored in files or databases is hard, if not impossible to model. Instead we assume that **Webservices** of various kinds may be used to access the results of simulations and other SimDB products.

Some of these will be standardised in the SimDAL specification, but custom services may also be introduced. The model allows one to describe the experiments and their results, which should allow users to discover results of interest, after which the web services can be called for actually accessing these.

3 Logical model overview

The actual data model that we propose here is a logical model in the sense of [33] based on the domain model described in 2.4. This reflects its origin as a model for the *Simulation Database* (SimDB), an application for searching for interesting simulations and related concepts. An actual SimDB specification will be created, but the logical model has value beyond that application. In particular the SimDAL protocol will use it in its queryData phase (SimTAP), albeit in a form customised for ease of querying for simulations of a limited set of experimental protocols.

SimDM is still implementation neutral, but it is fully detailed and represented in UML. It has a human readable HTML representation which contains the detailed description of all elements [5]. That document should be consulted for the details of the model.

Here we introduce the main concepts and motivate the main design decisions. Where possible we try to add a hyperlink from a concept’s name pointing into the HTML document the first time we use the name. The link will consist of a root

URL to the location of the HTML document, followed by a #<UTYPE> that identifies the description of the actual concept in the HTML document. This we feel is very much in the spirit of the use cases of UTYPEs. Later references to the concepts will in general not contain the link. Then class names will be capitalised. Abstract classes will be in italics. Names of packages, attributes, references or collections will be preceded by the class name where necessary, or it will be assumed to be clear from the concept what is intended.

3.1 Packages

UML Packages are subsets of classes and data types that are deemed to belong together. Whilst not essential to the model, we have used them to provide some level of modularity. Their main role is played in the XML schemas derived from the model. Each package has its own type-schema (see 4.2) which provides a somewhat finer level of reuse.

The diagram in Figure 2 shows the packages we use and their dependencies. This hierarchy is reflected in the UTYPEs, see section 4.1. The colours assigned to the packages correspond to the colours of classes in the diagrams in later sections. The subdivision in the one parent and three child packages follows the resource *class* hierarchy described next.

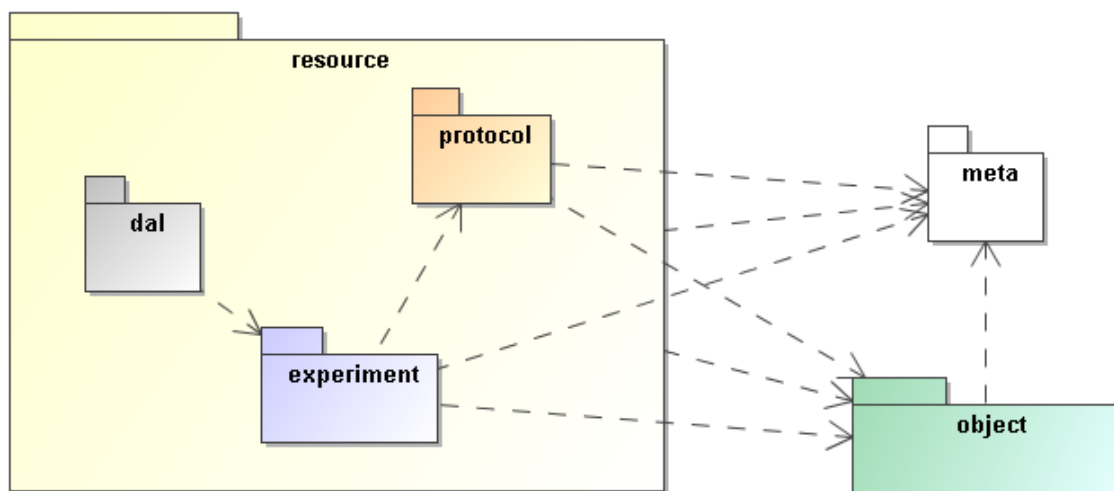


Figure 2: The packages of the SimDM and their relationships. These are related to each other through directed dependency links indicated by the dashed arrows.

3.2 Resource

The SimDM aims to describe simulations and related concepts. The current model does so with of the order of 40 separate object types, or classes. Most of these classes themselves represent parts of other classes. They group together properties or relationships used in the definition of their “parent”. The composition relation is used to represent these kinds of parent-child.

But among the classes in the model there are some that are not used like this. These classes represent concepts that can stand on their own, are not used to describe part of a larger concept. These we will call “root entity classes”. In the model they can be identified by the fact that neither they, nor any of their sub or base classes are part of another class, a child in a parent-child relation.

These are the classes that represent the model’s core concepts and their identification is a first important choice in the modelling effort. In the current model there actually two separate collections of classes that are root entities. The [Party](#) class represents an individual or organisation. It is used for indicating who/what wrote simulation codes or ran simulations. The main focus in this document is on the root entity classes in the [Resource](#) hierarchy, illustrated in Figure 3.

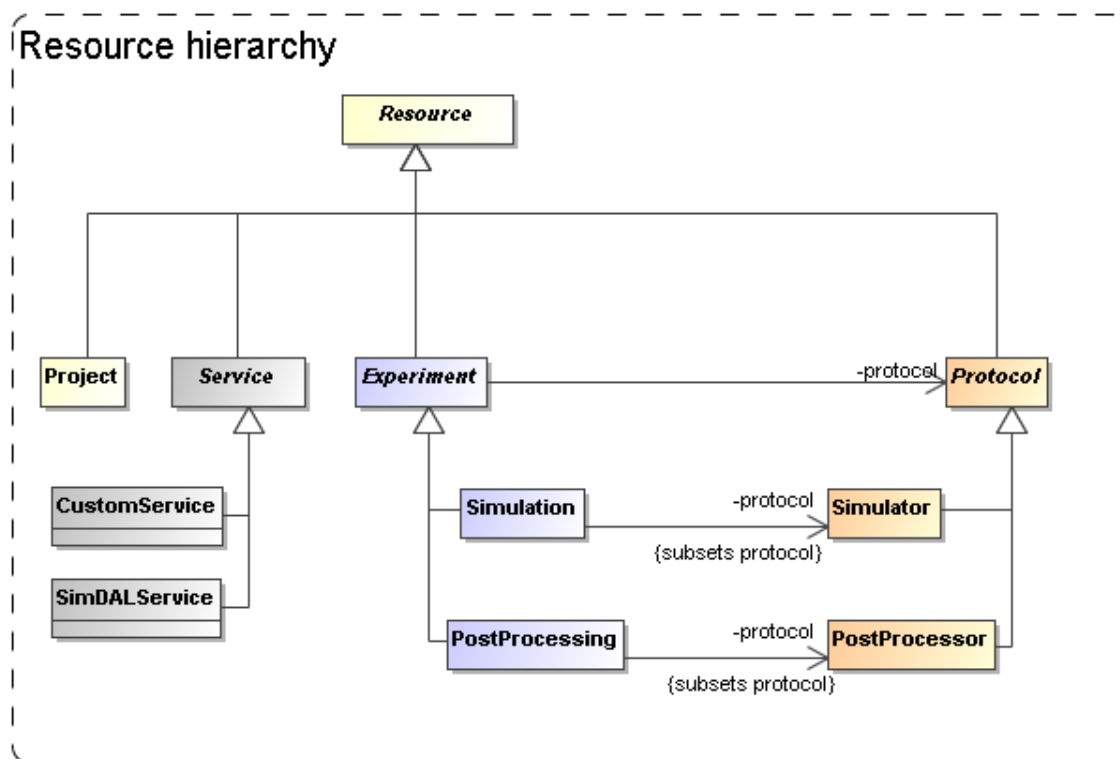


Figure 3: Root entity classes for SimDM.

From the top down we start with the ultimate root entity class, [Resource](#), which defines components common to all the main classes. The layer below it contains [Protocol](#)¹², [Experiment](#), [Service](#) and [Project](#). Protocol is the base class of the concrete classes [Simulator](#) and [PostProcessor](#). [Experiment](#) is the base class of [Simulation](#) and [PostProcessing](#). [Service](#) is the base class of [CustomService](#) and [SimDALService](#). Project has no subclasses and is concrete.

¹² The use of the name Protocol for the concept we introduce here has led to comments by some reviewers who feared confusion with the use of the same name in for example DAL protocols. In this document, the capitalised term Protocol will refer to the class in the model. When confusion with other usages might arise we may add use the phrase “(experimental) Protocol”.

Our choice for the root entities follows the domain model in concentrating on the scientific experiments as a whole. The *Experiment* class contains, amongst other components, classes representing the actual results (represented by the [OutputDataset](#) class) that people may wish to access. Those are *not* the core concepts in our model. This is in contrast for example to the spectrum data model [10] which focuses on the representation of the spectrum, and has the provenance and other metadata as sub-components.

One reason is that an experiment can exist without having (yet) produced any results, but to have results (as defined here) one always needs an experiment. This is a clear example of a parent-child dependency, where the child's life-cycle depends on that of the parent. The standard way to model such relationships is using a composition relation and that is how we have modelled it. More about the way we model results below in 3.7.

The separation between *Protocol* and *Experiment* is an important feature that we directly take over from the domain model. This design was already motivated in Section 2.4 and is related to the Measurement-Protocol pattern in [26]. That pattern says that when one does a *measurement* (of some property) it is important to remember the *protocol* by which the measurement was made ([26], p65). In [11] this was extended to experiments, which in general consist of large numbers of "measurements", all done in similar ways. Whereas the term measurement seems to be more applicable to observations, it is simple to generalise the concept a bit and apply it to the *calculation* of properties during a simulation. Actually this is similar to the CalculatedMeasurement in [26],

An important reason to keep this separation between *Experiment* and *Protocol* also in our logical model is to avoid having to redefine the parameters and other aspects of a simulation code each time a simulation is run.

The *Service* class did not appear in the original domain model in [11], but we introduced it in the model in 2.4 under the name *WebService*. In our model it represents a way to provide access to results of experiments. We could have tried modelling the way results are stored in files etc., but deemed it too complex to do so. This is in contrast for example to the spectrum data model, where we can model the data directly and even can predefine the representation of the data. There an access reference to the data files can be given to download a result. For simulations this is in general not possible. In many cases simulation codes have their particular proprietary formats, often storing single results over multiple files. Hence we merely allow users to describe services by which one can access results. Here bwe only make a separation between custom services and services following the SimDAL service specification thata is under construction in the IVOA in a collaboration between the theory interest group and the data access layer working group.

The *Project* class represents a scientific project, acknowledging that these in general use one or more experimental protocols to perform multiple experiments. This class is introduced to allow for example publishers to group simulations and post-processing runs that were produced with a common goal. It was inspired by a discussion on whether some of the SimDM/Resources could be registered as

Registry Resources as well¹³. Many of the simulations registered in a SimDB will not qualify for the same reasons that individual images do not qualify to be registered. Resources in an IVOA compatible registry are relatively coarse grained; correspond to archives full of images published through a SIAP service for example. A Project can be used to define such collections also in SimDM. And indeed one may wish to register such collections separately in a registry. In a data model one can use aggregations of the corresponding concepts to build such relations. In fact we have the single aggregation between Project and *Resource*, providing the user the freedom to include *Service*-s and even other Project-s.

The root of the hierarchy of entities is formed by the *Resource* class. This class is introduced as a convenience to hold on to information common to all its sub classes. Its name is obviously inspired by the Registry's Resource [12] and it also holds on to curation information. It "is not a" Registry Resource though in the strict OO modelling sense. For example it does not inherit all features of that class. But this is mainly because, as mentioned above, most SimDM Resources will not qualify as Registry Resources.

3.3 Physics, models and algorithms

An important characteristic of simulation codes is what physical systems and processes can be modelled and how these are represented in the program. The Simulator class represents computer codes that create numerical models of the world. Simulators do so by representing physical processes using numerical algorithms that act on model representations of real world objects. In our model, see Figure 4, physical processes are represented by the [Physics](#) class. It is contained in the Simulator class, not in the more general *Protocol*. In effect a simulation protocol is distinguished from other experimental protocols in that it models and implements physical processes.

Physical processes are implemented using particular [Algorithms](#). Algorithms are contained in *Protocol*, as also PostProcessors use them. In that case they implement the processing of existing results, and do not model physical processes. Examples of these are particular algorithms for extracting clusters from results of N-body simulations.

¹³ Thanks to Ray Plante for his contributions to this discussion.

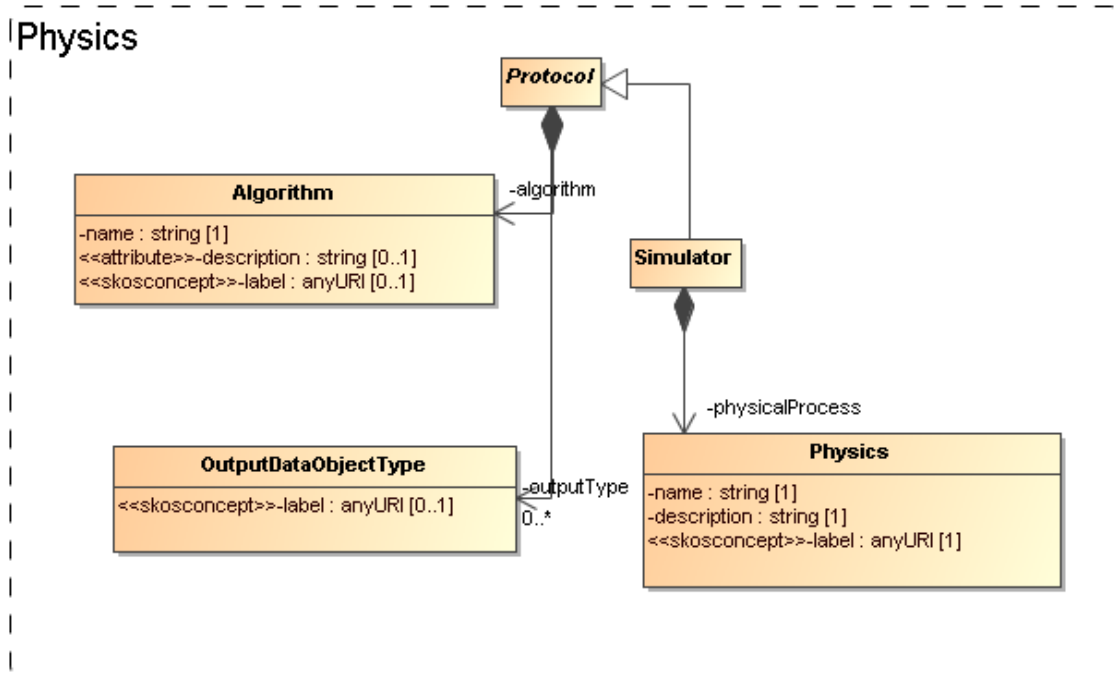


Figure 4: Modelling the representation of physical processes and objects.

Finally, experimental protocols need objects to represent the structure of the physical systems they model. For example, N-body simulations need particles that represent mass moving around. The model uses the [OutputDataObjectType](#) for this. This class allows one to define a hierarchy of data objects, from container objects like catalogues, data cubes or images down to the smallest objects such as particles or pixels

3.4 Parameters: definition and values

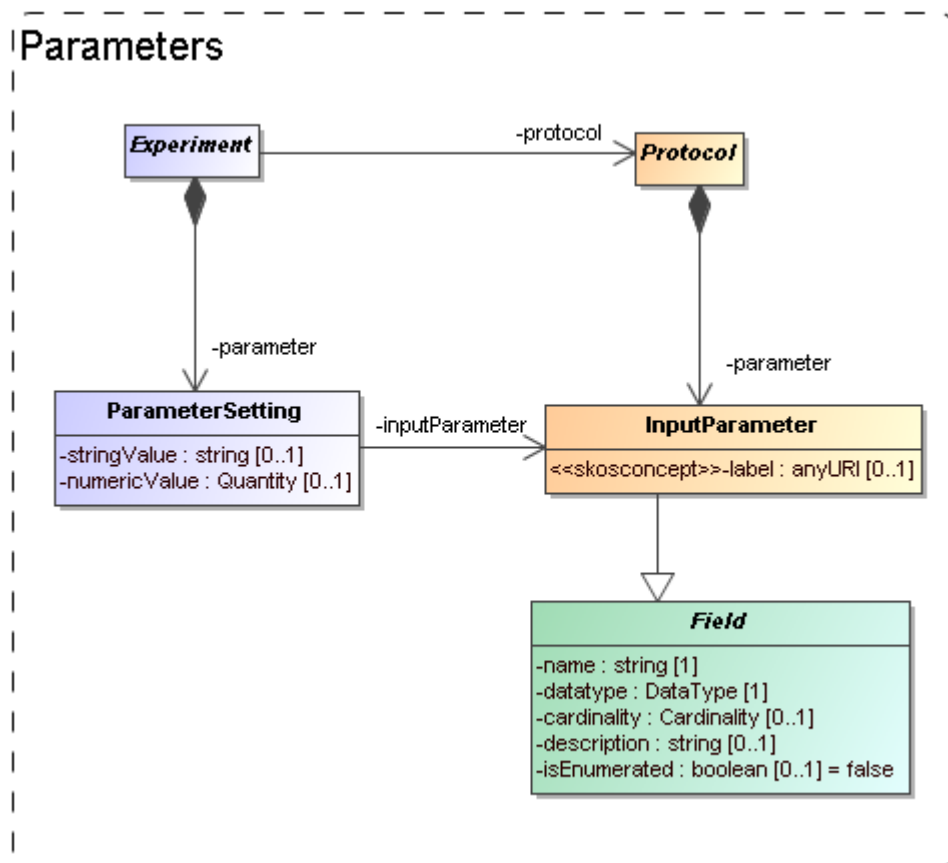


Figure 5: Modelling the parameters: definition under (*experimental*) Protocol, values under Experiment.

Software codes generally require some level of configuration before they are executed. In many cases this translates into a collection of parameters that must be given values. The parameters are defined by the code and we model this by an [InputParameter](#) class that is contained by *Protocol*. Assigning values to these parameters however is the responsibility of the experimenter and is explicitly modelled as a [ParameterSetting](#) class contained by *Experiment*.

Input parameters are defined by a attributes [name](#), [datatype](#), [label](#) and other properties familiar for example from the PARAM field in VOTable¹⁴. Most of these are inherited from the [Field](#) class, which will be discussed in Section 3.6 below.

Because the details of the parameter are defined on the InputParameter class, the ParameterSetting needs only a pointer (the [inputParameter](#) reference) to the appropriate input parameter and a value. A problem for this model though is what data type to assign to a possible value attribute. We have no knowledge in advance on the data type of the input parameter for which a value is set. This is only known at the instance level, not at the model level. We do not know whether

¹⁴ We generalize the ucd attribute on VOTable's PARAM and FIELD to a label attribute with stereotype <<skosconcept>>.

a certain parameter will be integer, or real, or maybe a string. Our current solution is to allow two different representations of a value, namely a [numericValue](#), of type real and a [stringValue](#) of type string. This issue and the usability problems it causes will need to be discussed and handled at the IVOA protocol level¹⁵. One approach is the SimTAP approach detailed in the SimDAL document [22].

3.5 Target: Goal of experiment

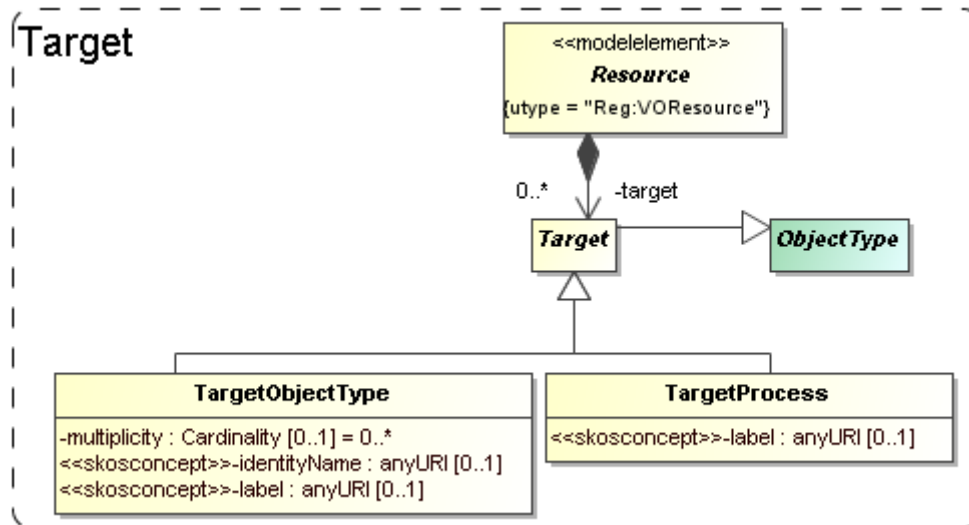


Figure 6: Modelling the goal, or target of a generic resource as objects and/or processes.

Generally the first piece of information that the scientists we polled were interested in regarding simulations was *what* was simulated. I.e. what type of object: a galaxy merger, a galaxy cluster, the large scale structure of the universe? This information in general says something about the goal that the scientists running the simulation had.

In certain cases the simulation code itself may completely prescribe the type of objects and physical processes that are modelled. As example take population synthesis models such as the Galaxev library¹⁶, producing spectra of galaxies.

But many simulation codes allow many different types of objects to be modelled, and even allow one to vary which processes are actually modelled. Also in many cases the actual object that is being simulated is not an intrinsic property of the simulation code, but is a derived property of the actual simulation. For example an N-Body code in general does not contain “galaxy particles”. But one can use it to follow the evolution of millions of low mass particles that are in a particular configuration that together model a galaxy. But it can also be a globular cluster, or a filament in the large scale structure.

To cover the concept of the target of an experiment or protocol, or the goal of a project, we add two classes, [TargetObjectType](#) and [TargetProcess](#). A

¹⁵ Statistical summary has the same problem, see 3.7

¹⁶ Bruzual and Charlot, 2003: <http://www.cida.ve/~bruzual/bc2003>

TargetObjectType represents an object, or a physical system in the real world, such as a galaxy, a star etc. TargetProcess represents a physical process such as gravitational clustering or turbulence. This recognises the fact that some simulations are run with the goal of investigating a process, rather than producing a model of a physical system.

Both these classes are subclasses of [Target](#), which itself is again a subclass of [ObjectType](#) defined in the next section. *Target* is contained by *Resource* so that by inheritance they are available to all sub classes. We do not model the *Target* objects in full possible detail. That we leave to future astronomical ontologies. We restrict ourselves to a [description](#) (inherited from *ObjectType*) and a semantic [label](#) attribute which identifies the intended concept using a standardised name from a SKOS vocabulary (see 5.2 below).

3.6 Object types: real and simulated

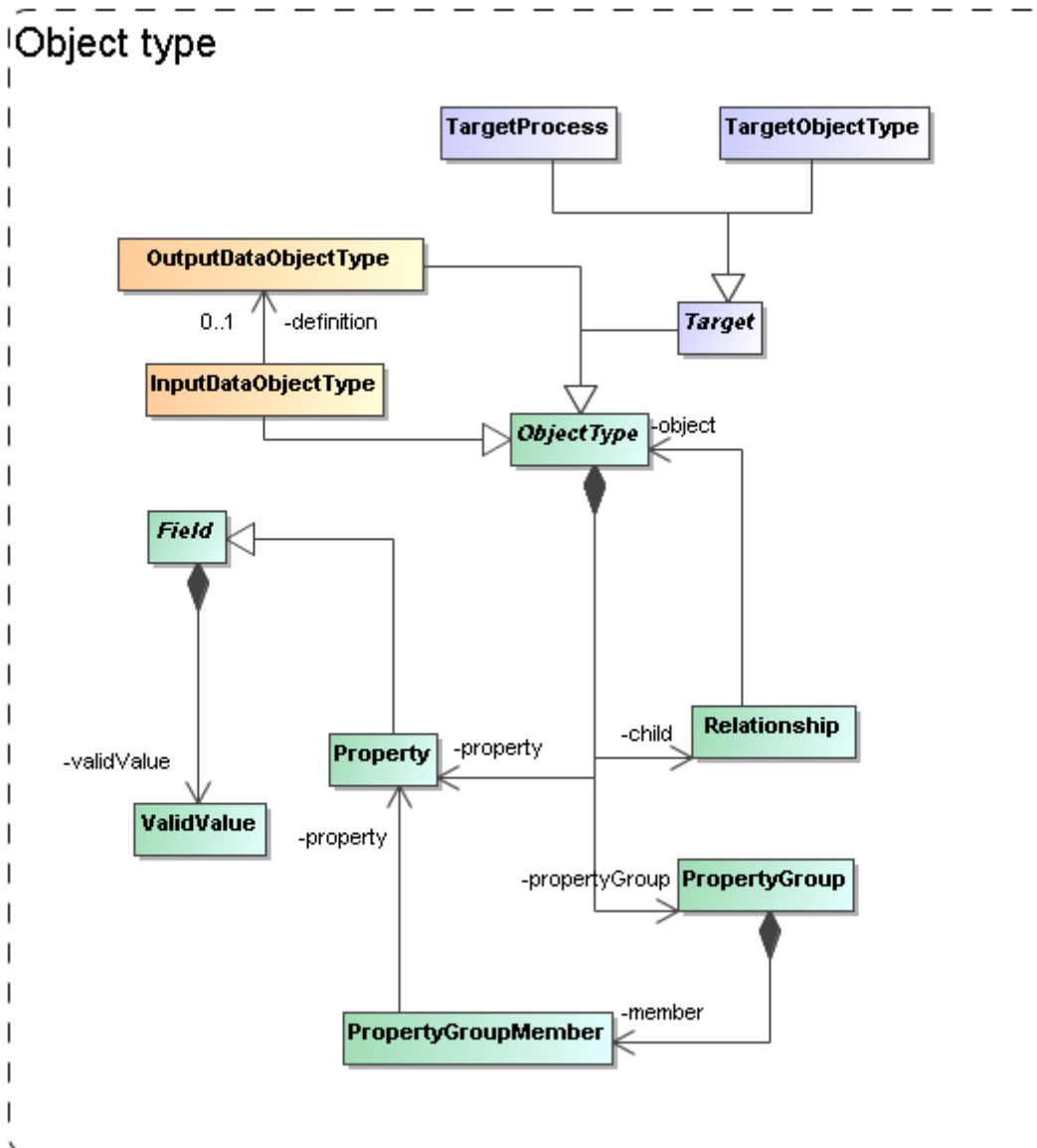


Figure 7: The model needs to describe the types of objects that are being simulated/used. We model this in quite some detail in a hierarchy of object types, with properties, grouping of properties and child objects corresponding to nested objects.

In a few places in the model we need to represent the fact that different simulation codes and other experimental protocols, or different experiments, need to describe the types of objects they use or produce. For example, the *Protocol* must be able to describe the building blocks of the model world it represents. We have encountered this need in the *OutputDataObjectType* in 3.3 and the *Target* (*ObjectType* and *Process*) in 3.5.

The building blocks required to describe a model world are entities with properties and relations. In SimDM we support this with a limited version of an object oriented meta-model

The core concept is the abstract [Object Type](#) class. An *ObjectType* contains a collection of [Property](#)-s that corresponds to the simple attributes used to describe an object. Property is a subclass of [Field](#) which defines its main attributes such as name, description and data type. Also a *Protocol's* InputParameter is a *Field*, similar to the way a VOTable's PARAM and FIELD share a common structure. Another similarity with the VOTable structure is the possibility to group Property-s in a [PropertyGroup](#).

To model (hierarchical) relations between different objects an *ObjectType* has a collection of [Relationship](#)-s, which can be used to define an aggregation or composition of other *ObjectType*-s. For example one can define that a simulation produces snapshots consisting of particles of different types, halo catalogues, containing halos, which themselves consist of particles, or images consisting of pixels.

3.7 Results: data sets and their statistical summary

We assume users of a Simulation Database will want to gain access to results of simulations and related experiments. This is the same as we assume of users of Simple Image Access or Simple Spectral Access services. For those services the user knows what to expect, a FITS image in one, a spectrum serialised according to the spectrum data model in the other.

Such expectations are not realistic for simulations though. The main problem is that we have no *a priori* knowledge about the contents of their results. Arguably somewhat simplistically one may claim that images and spectra contain pixels with known properties (space, wavelength, flux). Results of simulations, even when constrained to 3+1D simulations, can contain as their *fundamental constituents*: point particles, particles with size and structure, mesh cells of fixed or varying size, Voronoi cells¹⁷, structured halos, galaxies, radiation fields, galaxy merger trees etc. And any of these object types can come with any collection of properties: position, velocity, mass, temperature, chemical composition, entropy etc.

¹⁷ <http://www.mpa-garching.mpg.de/~volker/arepo/>

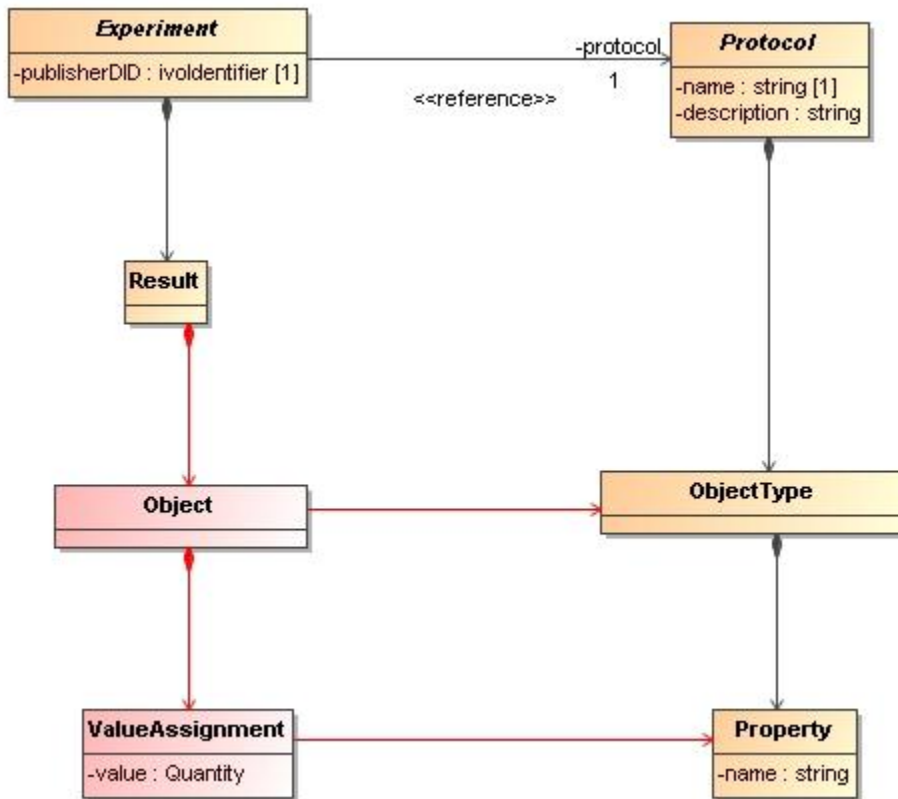


Figure 8: Domain model for results.

Precisely for this reason users will want to gain knowledge about the contents of simulation results to decide which simulations might be of interest to them. Hence the model must support description of the results explicitly. Figure 8 illustrates how this is achieved in the domain model: Experiments produce Results that consist of Objects (pixel, N-Body particle etc) of a particular ObjectType. The ObjectType defines the structure of Object as a collection of Properties (position, velocity, flux etc), and an Object, being an instance of the ObjectType, assigns values to these properties. Which ObjectTypes and Properties are available is defined by the (*experimental*) Protocol according to which the Experiment is run.

SimDM deviates from the domain model in that it does not include the Object and ValueAssignment classes. Including these would imply that positions, velocities etc for all particles in a simulation are added to the metadata description. For the purposes of SimDB, i.e. discovery of potentially interesting simulations, this would be overkill. It is certainly possible to include data as in the domain model, the spectrum data model is a case in point. But that model has a different purpose, namely providing a serialisation of spectra in a standard manner. Also, individual spectra are generally relatively small and have a well-defined structure.

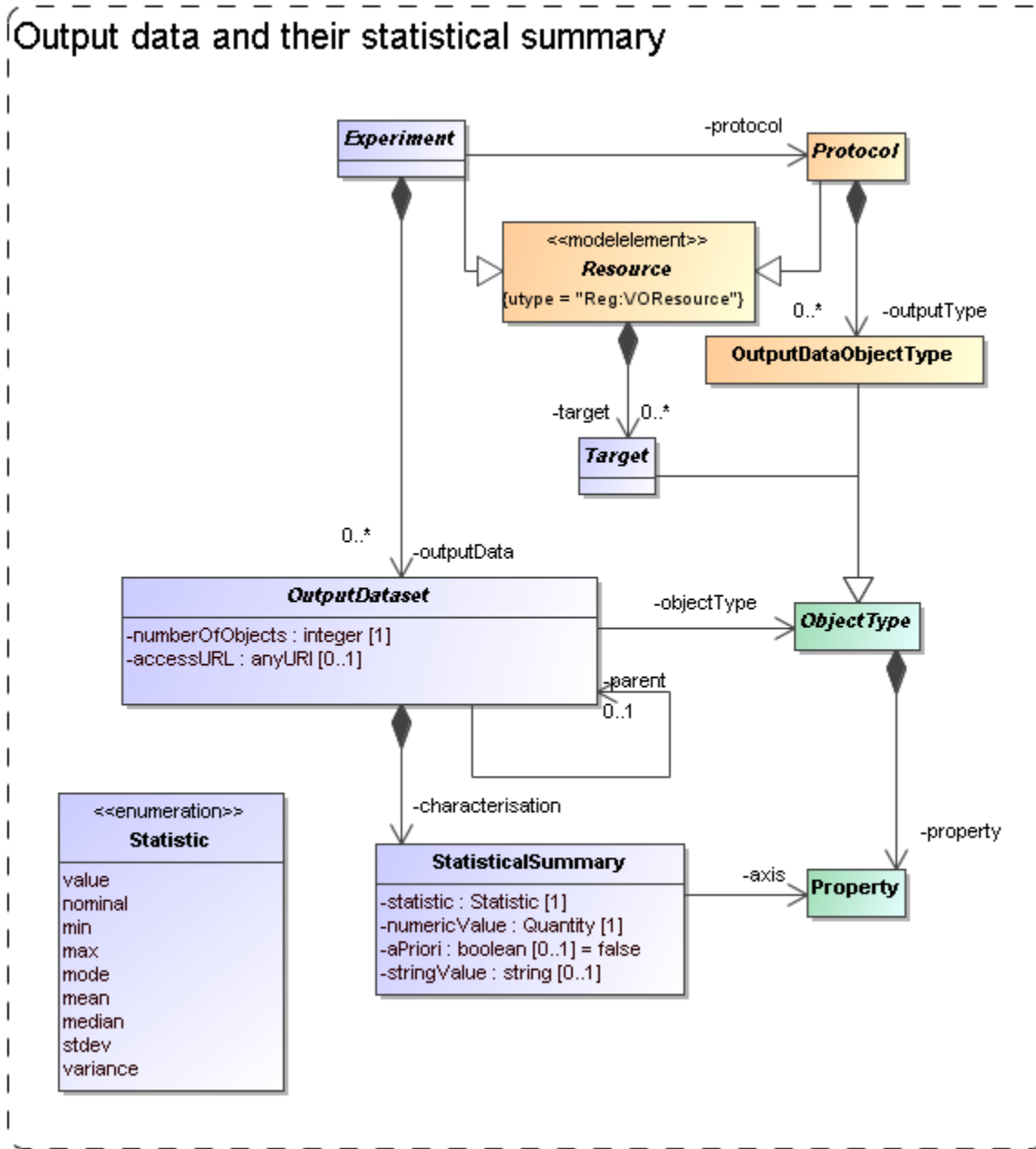


Figure 9: Modelling Output data sets and their contents.

How we model results and their contents in SimDM is shown in Figure 9. We model results as [OutputDataset](#)-s. An *OutputDataset* represents a *collection of objects* of a particular *ObjectType* produced by the *Protocol* during the *Experiment*, implemented using the [objectType](#) reference. In the typical case the referenced object will be an *OutputDataObjectType* defined for the given *Protocol*. But in principle it could also be for example a *TargetObjectType*, if a user wishes to describe properties of the astronomical objects that have been simulated.

The class hierarchy that users can define with *ObjectType* can be reflected in the definition of the *OutputDataset*. There the [parent](#) reference indicates the container object, corresponding to the container of the corresponding *Relationship* object in the *ObjectType* hierarchy¹⁸.

It is in certain cases useful to have some more quantitative information about the results of an experiment. For example, apart from the fact that a simulation has N-Body particles with properties position, velocity and mass, it might be of interest to know that the typical mass of the particles is 10^{10} solar masses.

We support this by allowing users to describe properties of the collections of objects in an *OutputDataset* using a class we call [StatisticalSummary](#). This reflects our belief that statistics is the appropriate way to introduce some quantitative aspects of these large collections of objects. *StatisticalSummary* is contained in *Product*. It assigns statistical values such as a mean or a min/max value to [Properties](#) of the [OutputDataset.objectType](#). Which statistic is used is described by the [statistic](#) attribute.

What data type to assign to a possible value attribute is again a problem, similar to that discussed for the value attribute in *ParameterSetting*. It is here solved in a similar way by introducing [numericalValue](#) and [stringValue](#) attributes. Whereas for a general SimDB a cumbersome solution such as this is necessary, the SimDAL protocol will allow for a more user friendly solution through its SimTAP query component (see [22]).

Extensions of this statistical summary to more detailed summaries such as histograms can be easily imagined, but have been left out for the model as they will have less relevance for discovery, which is the main use case for the model.

One further feature is important and is represented by the boolean [aPriori](#) attribute. This attribute describes whether the statistic that is used in the summary is an *a priori* or an *a posteriori* statistic. An *a posteriori* statistic is calculated using the results after they have been obtained during the running of the experiment. For example an *a posteriori* mean will likely correspond to the usual expression,

$$\frac{1}{N} \sum_i^N a_i ,$$

where the a_i are the values of some property.

In contrast *a priori* statistics characterise the possible values of the observables *before* the experiment is run. In certain cases *a priori* knowledge is available that restricts the possible values that certain properties may obtain in an experiment. An example is a lower bound set on the number of particles that a cluster must contain to be included in the result of a cluster extraction of an N-Body simulation. This can be indicated by a *StatisticalSummary* object with `statistic=min` and `aPriori=true`.

Knowledge about the *a priori* statistics is important in the interpretation of the results. In the previous example, when interpreting the mass multiplicity function

¹⁸ The assumption is that we generally do not need very complex relationships so that the simple parent reference is sufficient to infer which relationship is intended.

of a cluster catalogue extracted from an N-Body simulation, it is clearly important to know what the lower limit was on the mass of clusters.

In general *a priori* statistics are the result of, and may often be derived from the input parameters. However this derivation may not be obvious and will in general require intimate knowledge of the parameters of a (*experimental*) protocol. The *a priori* statistic may then facilitate the discovery of catalogues that should contain halos of a certain mass.

3.8 Data access services

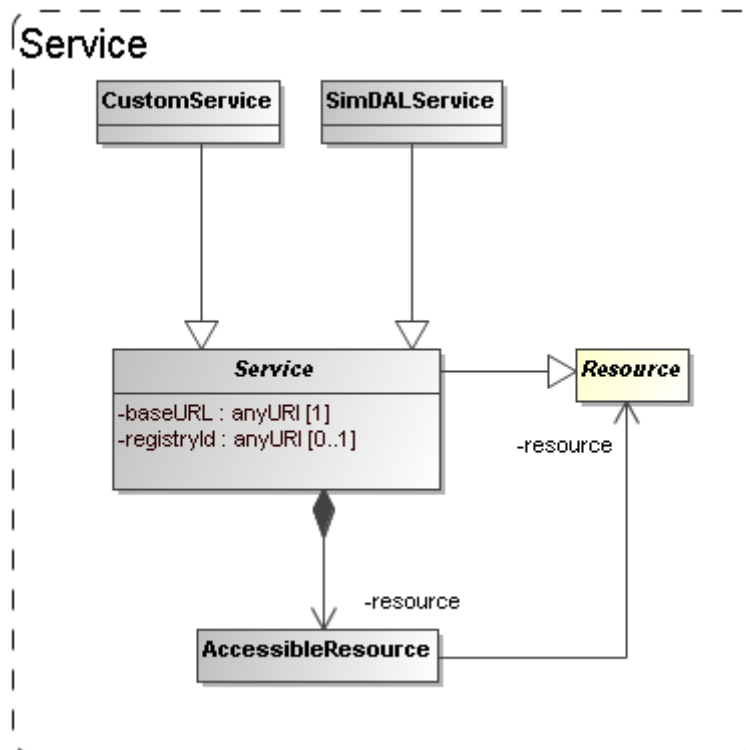


Figure 10: The model for (web) services giving access to SimDB resources.

The goal of the Simulation Database is to allow scientists to find simulations of possible interest. Once these are found the question is what can be done with them. Clearly knowledge of their existence will be useless if the researcher will not be able to somehow gain access to the results. The usual way this is done in the IVOA, for example in the simple image and spectral access protocols, is that the result of a discovery query contains an access URL that may be used to download the actual image or spectrum, where moreover the format of the returned resource, FITS, VOTable or XML document will be known beforehand.

It was perceived from the beginning of the SNAP project even that for the type of simulations that were supposed to be described a simple download would be unfeasible simply based on the size of many of the typical N-Body or AMR simulations. This assumption still holds and the SimDAL protocol is designed to

define special purpose services for retrieving parts of such simulations for example.

Also in the data model we want to indicate how the relation is between the results and services. This part may be used in the SimDB specification to allow users to register services and the Resources they give access to.

In the model the [Service](#) class, already introduced in 3.2, represents such access services (see Figure 10). This class can be explicitly linked to the Resources it gives access to through a collection of [AccessibleResource](#)-s. The class has concrete subclass [SimDALService](#). This represents services providing access to the results of a limited set of *Protocols* through the SimDAL protocol. The [CustomService](#) is introduced so that users can also publish non-standard services. This part of the model is still rather summarily treated and may need to be updated depending on developments in the SimDAL specification.

4 Serialisations

According to policies of the data modelling working group, first decided in Cambridge, 2003, a data model should be presented using a UML diagram, a corresponding XML schema and a list of UTYPEs. We have created these both using rules that derive the products directly from the XMI serialisation of the UML data model.

4.1 SimDM/UTYPE

The original goal of the data model presented here was to define the structure of a relational database supporting the Simulation Database *service* specification. A first draft of a note proposing that spec can be found in [21]. SimDB will use TAP [9] to define the IVOA protocol for querying this database using ADQL. The results of such queries will be tabular and serialised as VOTables. Such a VOTable will contain a filtered subset of the information in the database, but in general in a different form compared to the structure of the data model. To indicate the meaning of data elements in such a VOTable, the IVOA has invented the concept of UTYPEs.

A UTYPE is a “pointer into a data model”¹⁹. The VOTable XML schema implements this concept as attributes on various elements, e.g. FIELD and TABLE and many other elements. The value of such a UTYPE attribute should identify an element in a data model that is represented by the element itself. For example a table might point to a class definition in a data model, and a column (FIELD) to an attribute.

It has become common practice to provide for an IVOA data model a list of UTYPEs. The Spectrum data model (see [10]) was the first to add explicit UTYPE-s for each of the attributes in its model and the Characterisation data model [15] has followed that example. We follow these examples by assigning UTYPE-s explicitly to all elements in the model.

¹⁹ See 5.3.3 for our position on the discussion that is still going on regarding UTYPEs.

Our goal was not to have to make this a separate effort, but if possible to generate the list of UTYPEs directly from the model. Our goal was to assign UTYPEs to all identifiable elements in our model and these should be unique.

To this end we define a set of production rules phrased using the special names in our UML profile. We have made a guess as to what the format for UTYPEs will be. In the previous data models a UTYPE existed of a word consisting of dot-separated “atoms”, similar to UCDs, but without the “;”. We use a slightly different format to make the distinction between different syntactic elements from the profile somewhat clearer and also to guarantee uniqueness of each UTYPE within the data model context. Once (if?) a format is settled on within the IVOA we will easily be able to adjust our definitions.

The important point we want to make is that it is possible to define simple rules that can automatically produce *unique* UTYPE-like words for all elements of a data model, i.e. the only discussion that may be required is on the rules for doing so IF a fixed format is preferred (see Norman Gray’s ideas²⁰ on why this might not be necessary).

The following BNF-like expressions define the particular rules we have used for deriving the UTYPEs from the UML model:

```
utype          :=      [model-utype | package-utype | class-utype |
                        attribute-utype | collection-utype |
                        reference-utype | container-utype

model-utype    :=      <model-name>
package-utype  :=      model-utype ":" package-hierarchy
package-hierarchy := <package-name> ["/" <package-name>]*
class-utype    :=      package-utype "/" <class-name>
attribute-utype := class-utype "." attribute
attribute      :=      [primitive-attr | struct-attr]
primitive-attr := <attribute-name>
struct-attr    :=      <attribute-name> "." attribute
collection-utype := class-utype "." <collection-name>
reference-utype := class-utype "." <reference-name>
container-utype := class-utype "." "CONTAINER"
identifier-utype := class-utype "." "ID"
```

For the SimDM these rules produce a list of UTYPEs for the model. For each model element we provide the UTYPE in the HTML documentation in [5] and we provide a complete list at the end of that document²¹. Note also that a URL of the type

`<URL-to-HTML-doc>#<utype>`

will link one directly to the documentation for the corresponding data model element. This is in conformance with a suggestion made by Norman Gray²⁰.

²⁰ <http://nxg.me.uk/note/2009/utype-proposals/>

²¹ <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/html/SimDM.html#utypes>

When representing components of the data model in a VOTable (for example), these UTYPEs SHOULD be used, in particular when the VOTable contains results of ADQL queries to a SimDB/TAP implementation (see *SimDB Services*). Alternative views and representations of the SimDM, for example in SimDAP, SHOULD use these UTYPEs to refer to elements in the model.

4.2 XML

A specification for an IVOA data model should (must?²²) contain an XML schema that defines how to serialise data model instances as XML documents. Similar to the case of UTYPEs we did not want to make the design of these schemas a separate effort; instead we want to derive the schema from the model. To do so we have defined rules for relating XML Schema constructs to our UML model. These rules are a completion of those described in [36]. It is based also on a view of what such schemas should look like, restricting the possible set of constructs to be used in schemas representing data models. These design rules have earlier been discussed with and accepted by the Registry and VOTable working groups.

We give here only a short description of these rules. First of all we define two different types of schemas. First we define “type schemas”, XSD documents containing only type definitions. For each object type(class) and value type we generate a corresponding complexType or simpleType. Attributes map to elements of a corresponding data type (simple or complex), collections to elements of a type corresponding to the class. References are harder to represent and will be discussed below.

We next generate a “document schema” containing root elements. The elements in the document schema define the valid XML documents one can write and we choose only “root-entity classes” for their type. That is, only classes at the root of collection trees can be represented as a document. Fragments of these are not allowed. For example, only a complete Simulator or Simulation can be represented in a document, not only a single result, or parameter setting.

Note that this is a choice made for the Simulation Database service specification. The document schema depends on the type schemas through XML schema import declarations. This separation allows flexible usage of the type schemas, for example other services might make a different choice from the types to serve as valid root elements.

The root schema for the SimDM/XSD representation can be found [here](#)²³. The type schemas and a predefined base schema can be found in the same directory and subdirectories of it. We refer to the *SimDB Services* document for more details on the XML schema serialisation and their use in the SimDB service protocol.

²² See “Rules” on <http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/HowToParticipate>. This “decision” was made in the Cambridge 2003 interoperability meeting together with the requirement that data models must be specified in UML.

²³

http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/xsd/SimDM_root.xsd

Only the mapping of references deserves special attention. Our choice of mapping from UML to XSD elements and our definition of root elements imply that many references must be able to link between different XML documents. For example the (*experimental*) protocol reference²⁴ in an XML document describing an Experiment must be able to identify a (*experimental*) Protocol that is defined in a different XML document. To do this identification we assume we must rely on an agent that can interpret a serialisation of a reference and use it to look up a corresponding document. Therefore we map references to elements of a particular complexType that we define in a base schema²⁵. That same schema defines a type to be used for representing identifiers of objects and the reference serialisation must be able to reproduce such an identifier. Further technical details of this mapping will be described in the appropriate service definition document.

5 Dependencies on other IVOA efforts

IVOA documents are assumed to specify dependencies on other IVOA efforts. We have from the beginning realised that the SimDB effort touches upon various other specifications and general efforts of other working groups [21]. Here we discuss these relations as far as they pertain to the Simulation data model.

5.1 Registry

The correspondence between the full Simulation Database specification and the IVOA Registry will be discussed in the *SimDB Service* note [21]. Here we will address the relation between the SimDM and the Registry Data Model as defined in [13].

²⁴ UTYPE: SimDM:/resource/experiment/Experiment.protocol or <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/html/SimDM.html#SimDM:/resource/experiment/Experiment.protocol>

²⁵ <http://vo-urp.googlecode.com/svn/trunk/xsd/base.xsd>

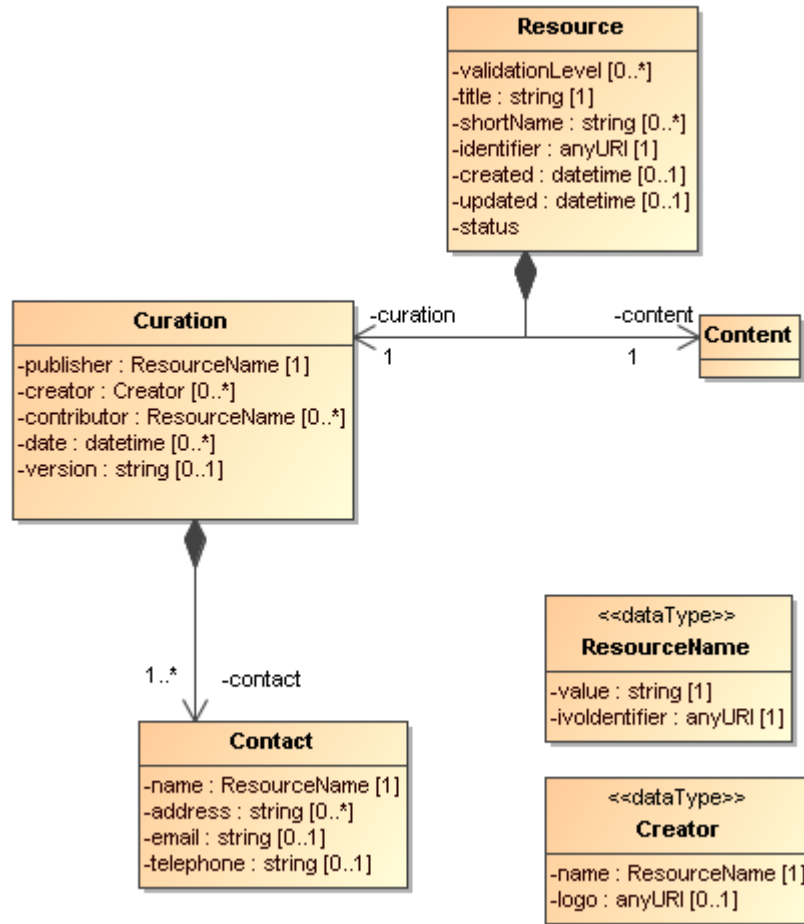


Figure 11: UML rendering of the Resource complexType from [13].

In Figure 11 we present a UML rendering of the *Resource* complexType as inferred from the Resource Registry VOResource XML Schema [13]. Comparing that model to SimDM/Resource we can see that these two models for Resource are related, but not identical. In data modelling terms, it is not true that a SimDM/Resource *is a* Registry/Resource (or *vice versa*). *Curation* is modelled differently and arguably with less detail in SimDM, but the main difference is in the *Content*. SimDM provides a very detailed and specialised model for the *Content* of Simulations and related resources, by modelling provenance, motivation and results characterisation. This higher level of detail gives rise to a higher level of granularity in the types of resources stored in a SimDB, which in general will be too fine grained for registration in a Registry. This is similar to the case of a single image, which is not a Registry/Resource, whereas a SIAP-compatible *service*, providing access to many images, is.

A SimDB service itself will have to be registered, i.e. a SimDB service *is a* Registry/Resource. In discussion with Ray Plante (IVOA Interoperability meeting May 2007, Beijing) on this issue it was proposed that some part of the contents could also be registered in a Registry directly, i.e. we should be able to identify Registry/Resource-s in SimDB. Considerations to decide on how to make this identification would be for example that all data products resulting from a well-

defined (and published) scientific project could qualify. To represent such a possibility for now we have introduced another subclass of SimDM/Resource: SimDM/Project. This is not much more than an annotated aggregation of other SimDM/Resources, with some additional attributes describing the motivation etc. The metadata of a SimDM/Project is not the same as that of a Registry/Resource, however we propose that we should be able to define a transformation (possibly implemented again in XSLT) to transform a SimDM/Project and produce a Registry/XML representation.

5.2 Semantics: Use of SKOS Concepts

In the SimDM, observables, object types, properties, parameters that play a role in a given simulation have to be defined explicitly, for the world of simulations is too large to define all possibilities explicitly in the model itself. This in contrast for example to the spectrum data model [10] where we know that a flux is determined for a wavelength interval, or a model for images where a flux is determined for a spatial pixel. In principle the publisher of a SimDM/Resource has all freedom to name and describe these entities. For other users to understand the meaning of them, we have where appropriate, added an attribute corresponding to a semantic label. This is similar to the situation in VOTable, where FIELD-s can be given a UCD (or UTYPE) that allows users to understand the meaning of a column in the table.

In SimDM we need to generalise this concept as UCDs are not sufficient for our purpose. For example target object types are not covered by the list of UCDs and the same for other elements in our model. The Semantics WG has specified that such vocabularies should follow the SKOS specification [24]. They have also defined a number of such semantic vocabularies in the SKOS format, for example of astronomical objects. We try to anticipate their results by introducing a special type of attribute in our UML profile that corresponds to a concept in a given ontology.

Technically, in the UML profile we have defined a stereotype <<skosconcept>> that can be assigned to an attribute in the UML model. Attributes with this stereotype must define a value for the tag "broadestSKOSConcept".

The intent of this is as follows (thanks to Norman Gray for providing the original text with this formal definition):

<<skosconcept>> attributes take a skos:Concept as their value. In each case, the value is given as a single skos:Concept: such attributes may take any skos:Concept which is a narrower concept than this single typing concept. To be precise, for a typing concept T, any concept c is a valid value for this property, if either:

c skos:broaderTransitive T

or if there exists a concept X such that

c skos:broaderTransitive X. X skos:broadMatch T

This just means that, if *c* is in the same vocabulary as *T*, then it's connected by a chain of any number of `skos:broader`, and if it's in a different vocabulary, then there is some *X* which is in the same vocabulary as *c*, with a cross-vocabulary link between *X* and *T*.

In several cases -- particularly those vocabularies which have been created for SimDM -- there will be a single top concept which everything is narrower than. In other vocabularies -- such as the AstroObject in the thesaurus version of the ontology of object types -- the natural typing concept is not a top concept, or is not the only top concept. This definition also does indicate that it's legitimate for concept *c* to come from a different vocabulary from *T*: the fact that *c* has been declared to be narrower than *T*, either implicitly or explicitly, is to be taken to be the expression of the vocabulary designer's intention that this be a legitimate value for this property.

5.3 Data Model

5.3.1 UML Profile

The data model proposed in this document is fully defined in all detail through a UML model. UML is a large language and we have consciously restricted ourselves to a subset of the possible modelling elements. We have also added a few modelling elements using the extension mechanisms UML provides through stereotypes, tags and predefined data types. This combination of restriction and extensions is referred to as a UML Profile. The details of our profile are described in a separate document **Erreur ! Source du renvoi introuvable.**, added as an Appendix to the current WD.

One reason to put so much emphasis on the UML model is that it allows us to derive various products of this specification automatically. To this end we use the modelling frame work under development in the VO-URP project²⁶, which is a spin-off of the SimDB effort. Using XSLT scripts developed in VO-URP we can generate HTML documentation (including UTYPE lists) [5], XML schema definitions [6] etc directly from the XMI representation of the UML model.

5.3.2 Characterisation data model

As described in section 3.7, the model allows one to characterise the results of experiments statistically using the `StatisticalSummary` class. This part of the model addresses similar problems for simulations as does the Characterisation Data Model for observations. We have not followed that model in detail, but have tried to incorporate its main ideas, giving a new interpretation to some of these²⁷. We believe the best way to reconcile the two approaches is to see both as specialisations of a more abstract model defining statistical characterisations of

²⁶ <http://code.google.com/p/vo-urp/>

²⁷ This follows ideas presented in China 2007, see

<http://www.ivoa.net/internal/IVOA/InterOpMay2007DataModel/CharacterisationInTheDomain.ppt>

data products. A proposal for such a “domain model for characterisation was given in [31].

5.3.3 UTYPE

Section 4.1 describes how we generate UTYPEs for the different elements in our data model. The rules we use to do so have been subsumed in a draft for a Note on UTYPE-s by [16]. One problem we have with that Note is that the concepts used in the grammar, and that are direct reflections of syntactic modelling elements in our UML profile, have not been defined. For models defined with different UML syntax the grammar does not help.

Some have argued against any semantic meaning to a UTYPE string. It should not be necessary to parse it to find out what its meaning is. Instead one should be able to follow it, but could/should be opaque. It should simply be assigned to the modelling elements. In that case the only requirement would be that a unique list of strings is created and that

Our assumption has been that a UTYPE should allow one to uniquely identify a concept in a data model. We do not assume that our particular form to do so need to be taken over. But, as we describe in 4.1, *if* one wants to simply derive a list of unique strings to be associated to concepts that play a role in data models designed with our UML profile, these rules may help. Clearly if the syntax were to change we can accept that.

The effort on understanding what UTYPEs really are, how they are to be used, or defined is in our opinion not completed. But we feel that our approach is compatible with any possible interpretation, and sufficiently flexible to proposed changes in precise syntax, were they required.

6 References

6.1 Accompanying documents

- [1] This document, at web address
<http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/PR-SimulationDataModel-v.1.00-20111019.doc>
- [2] SimDM UML diagram obtained from MagicDraw :
http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/uml/SimDM_DM.xml
- [3] A PNG representation of the main diagram, ‘all’, in the model, extracted from MagicDraw in
http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/uml/SimDM_DM.png
- [4] “Intermediate representation” of the model. An XML document containing all relevant information from the model in a more readable format than XMI. This document is generated from the XMI and is itself the source of all other generated products.

- http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/uml/SimDM_INTERMEDIATE.xml²⁸
- [5] HTML representation of the SimDM in
<http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/html/SimDM.html>
- [6] XML schema documents derived from the data model and defining the representation of data model instances in XML. Divided over various documents. The “element schema” document defining all root elements can be found here: http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/xsd/SimDM_root.xsd . All type schemas can be found in the same folder, <http://volute.googlecode.com/svn/trunk/projects/theory/snapdm/specification/xsd> and sub-folders of it.

6.2 Relevant IVOA documents

- [7] Claudio Gheller, Rick Wagner et al, *Simulation Data Access Protocol (SimDAP)*, <http://code.google.com/p/volute/source/browse/trunk/projects/theory/snap/SimDAP.html>
- [8] Bob Hanisch, *IVOA Document Standards*, <http://www.ivoa.net/Documents/latest/DocStd.html>
- [9] Pat Dowler, Guy Rixon, Doug Tody, *Table Access Protocol* <http://ivoa.net/Documents/TAP/>
- [10] Jonathan McDowell et al (2007) *IVOA Spectral Data Model* <http://www.ivoa.net/Documents/latest/SpectrumDM.html>
- [11] Gerard Lemson, Pat Dowler, A.J. Banday, 2004 *A Unified Domain Model for Astronomy* http://www.aspbbooks.org/a/volumes/article_details/?paper_id=861 see also <http://www.ivoa.net/internal/IVOA/IvoaDataModel/DomainModelv0.9.1.doc>
- [12] Bob Hanisch et al, *Resource metadata for the virtual observatory* <http://www.ivoa.net/Documents/latest/RM.html>
- [13] Ray Plante et al 2008, *VOResource : an XML Encoding Schema for Resource Metadata* <http://www.ivoa.net/Documents/REC/ReR/VOResource-20080222.html>
- [14] Carlos Rodrigo et al, *S3 : proposal for a simple protocol to handle theoretical data (microsimulations)* <http://www.ivoa.net/Documents/latest/S3TheoreticalData.html>
- [15] Mireille Louys et al (2008) *Data Model for Astronomical Data Set Characterisation Version* <http://www.ivoa.net/Documents/latest/CharacterisationDM.html>
- [16] Mireille Louys et al (2009) *Utype : A data model field name convention Version 0.3* <http://www.ivoa.net/internal/IVOA/Utypes/WD-Utypes-0.3-20090522.pdf>
- [17] Paul Harrison et al, *Simple Image Access specification Version 1.0* <http://www.ivoa.net/Documents/SIA/>
- [18] Doug Tody et al, *Simple Spectral Access specification version 1.04* <http://www.ivoa.net/Documents/latest/SSA.html>
- [19] *Theoretical Spectral Access Protocol* <http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/IVOATheryTSAP>

²⁸ The XML schema file defining the structure of the representation in [4] can be found in the VO-URP project: <http://vo-urp.googlecode.com/svn/trunk/xsd/intermediateModel.xsd>

- [20] Theory in the VO
G. Lemson and J. Colberg (2003)
<http://www.ivoa.net/pub/papers/TheoryInTheVO.pdf>
- [21] Gerard Lemson et al (2008) *Proposal for a Simulation Database Standard*,
IVOA Note 11 July 2008
<http://www.ivoa.net/Documents/latest/SimDBTrack.html>
- [22] Franck Le Petit et al (2011) *Simulation Data Access Layer* (in preparation)
https://volute.googlecode.com/svn/trunk/projects/theory/simdal/Simdal_draft.doc
- [23] Inaki Ortiz et al (2008) *IVOA Astronomical Data Query Language*
<http://www.ivoa.net/Documents/latest/ADQL.html>
- [24] Sébastien Derriere et al (2009) *Vocabularies in the Virtual Observatory*
<http://www.ivoa.net/Documents/latest/Vocabularies.html>

6.3 Other sources

- [25] Santi Cassisi et al (2008) *Framework for the inclusion of theory data and services in the VObs*
http://cds.u-strasbg.fr/wikiDCA/pub/EuroVODCA/Deliverables/EuroVO-DCA_D11_MPG_Final.pdf
- [26] Martin Fowler (1997) *Analysis Patterns*
Addison Wesley Longman, Inc
- [27] Terry Halpin (2001) *Information Modeling and Relational Databases: From Conceptual Analysis to Logical Design*
Morgan Kaufmann Publishers
- [28] XML schema, <http://www.w3.org/XML/Schema>
- [29] *MOF 2.0/XML Mapping, V2.1.1*
<http://www.omg.org/spec/XML/2.1/PDF>
- [30] *OMG Unified Modeling Language (OMG UML), Infrastructure Version 2.2*
<http://www.omg.org/docs/formal/09-02-04.pdf>
- [31] Gerard Lemson (2007) *Characterisation in the domain*
Presentation at IVOA interoperability meeting Beijing 2007.
<http://www.ivoa.net/internal/IVOA/InterOpMay2007DataModel/CharacterisationInTheDomain.ppt>
- [32] http://en.wikipedia.org/wiki/Conceptual_data_model
- [33] http://en.wikipedia.org/wiki/Logical_data_model
- [34] http://en.wikipedia.org/wiki/Physical_data_model
- [35] Jim Gray et al (2002) *Data Mining the SDSS SkyServer Database*
<http://www.sdss.jhu.edu/ScienceArchive/pubs/msr-tr-2002-01.pdf>
- [36] Gerard Lemson (2004) *Model Based Schema*
PPT presented during Registry video conference 2004-05-13
http://www.gvo.org/www/uploads/Documentation/Registry_XSD_videocon20040513-14.ppt

Appendix A History

Numerical computer simulations form an increasingly important component of astrophysical research. Such simulations are used to model astrophysical processes whose complexity precludes an analytical treatment. The subject of these simulations includes every possible astrophysical phenomenon, from the structure of stellar atmospheres, the formation of solar systems, the structure of galaxies and the description of their constituents, to the formation of the largest structures in the universe.

The simulations often result in predictions that can be compared to observations, but in general are much richer, including “observables” that can only be derived by indirect means from observations. These results can be very large, rivalling and often exceeding in size the largest observational catalogues. But they can also be relatively small, consisting of individual spectra of say a white dwarf, though often in collections resulting from parameter studies.

The design and execution of these simulations has become a specialised field of astrophysics, and is these days often performed in large collaborations. And while it is still true that their results are studied by these groups only, more and more of these theoretical data are being published online (see for instance the Appendix B of [25]).

Apart from limited support for publishing theoretical spectra in SSAP, there is as yet no IVOA standard dealing with the publication of simulations and their results. In earlier documents we have described the issues for defining such standards compared to the arguably simpler case of observational data sets (see for example [20] and [25]).

The proposal for a standard way of publishing simulations was formulated during a workshop in Cambridge, February 2006. The original idea was to create an analogue of the simple image access protocol (SIAP, [17]) for N-Body simulations: SNAP, the *Simple Numerical Access Protocol*. During the following interoperability meeting in Victoria, May 2006, the scope was expanded to include other types of simulation algorithms, and rephrased to something like “simulations that reproduce 3+1dimensional space time”. It was felt furthermore that not only simulations themselves should be included, but also certain types of post-processing such as cluster finders, as long as their results are still aimed at producing a description of 3D space at one or more points in time. Over time requests have come in to generalise this scope even more, basically to enable any type of astrophysical simulation to be handled.

An important change that was decided in Victoria 2006 was that instead of the SIA protocol, the newer simple spectral access protocol (SSAP, [18]) should be followed as an example. This protocol’s main difference with respect to SIAP was the explicit data model that was created for spectra and was used as motivation for the queryData metadata and the getData data format. Hence SNAP from the

beginning had a double focus on a data model plus related query protocol on the one hand, and a data access and delivery specification on the other hand. Shortly before the Trieste interop in the spring of 2008, it was decided to split SNAP up along these lines in two separate specifications: a specification for a *Simulation Database* (SimDB) which would support searching for interesting simulations and services providing access to them, and a *Simulation Data Access Protocol* (SimDAP) providing a specification for accessing simulation results.

SimDB on its own is still a rather complex specification. It has overlap with the efforts and results of many working groups, Data Model (DM), Registry, Data Access Layer (DAL), Semantics as well as being an integral part of the Theory Interest Group (TIG). This issue has been discussed in the Baltimore and Strasbourg interops, as it causes a potential problem for the standardisation process: an interest group cannot promote a document to a standard, but which a working group (WG) could do so. It was decided in Baltimore to postpone that decision by creating a focus group led by the TIG and with participation from the various WGs.

The current document is the result of a split in original Note that was written for SimDB. Such a split was proposed to simplify the standardisation process and after some refactoring was performed mid-2009. This current document is the first of these and deals exclusively with the data model (SimDM) and consequently has a natural place in the DM WG. The second document deals with the use of the data model for defining the model for a relational database and its related TAP query implementation as well as a service interface for uploading simulation descriptions to this database. It is not yet clear whether it can find a place in a single WG.

A parallel effort has been the proposal for a simpler access standard for small scale simulation, the Simple Self-describing Service protocol (S3, [14]). This was a result of an investigation started in the Cambridge 2007 interoperability meeting whether “micro-physics” simulations as they are sometimes called require special attention. For some time this was covered by SSA, at least as far as theory spectra were concerned. S3 is actually a direct reworking of an older Theoretical Spectral Access Protocol [19].

There were questions in the TIG whether S3 might be incorporated in SimDB and/or SimDAP. In the interoperability meeting in Victoria 2010 the decision was made that indeed this should be possible. The SimDM was shown to be able to incorporate the metadata for S3-like services, and indeed proposes extensions of that. It was decided that the S3 protocol should be merged with the SimDAP protocol, which from then on will be known by the name Simulation Data Access Layer (SimDAL). The implementation note addresses this question from a formal point of view, namely by defining *how* S3-like services can be described by the data model.

Appendix B UML Profile²⁹

The Simulation Data Model uses UML as the language for its specification. This is in accordance with decisions of the IVOA data model working group. One advantage of UML is that it is implementation neutral. It is a graphical language, consisting of “boxes and lines” that is very suitable for whiteboard discussions but allows one to model the concepts and relations in a static data model. It is also rich enough to allow one to describe all important data elements and relations.

In fact, UML is almost too rich. It is easy to become overwhelmed by the large number of possible syntactic elements to choose from for modelling a particular structure. Luckily UML allows one to formally define a subset of its language where one restricts oneself to a subset of the syntactic elements. Such a subset is called a *UML Profile*. Apart from creating a more restricted language, a Profile also allows one to assign new meanings to existing elements by defining *stereotypes* with associated properties (*tag definitions*). It is also possible to predefine classes and data types (see below) that can be reused by the data modeller.

In our modelling effort we have defined an initial implementation of a UML profile as created by MagicDraw. The profile³⁰ is contained in the UML file containing the SimDM data model. Here we give a list of the main elements that we use and give a short motivation for their inclusion in the language. It is our opinion that the DM working group should be ultimately responsible for a profile such as this, as it gives the possibility of defining a domain specific language for all IVOA data modelling efforts, thus giving some uniformity to those disparate efforts.

B.1 Element

All elements mentioned below are specialisations of UML Element.

Stereotypes

- **<<modelelement>>**: This stereotype can be assigned to any UML *Element* and is used to define the **utype** tag on.
Tags:
 - **utype [string]**: this holds the actual UTYPE that points to the other modelling element that is represented here.

B.2 Model

²⁹ This Appendix could eventually be replaced by a Note on UML profiles if the DM WG is interested in organising such an effort. As this does currently not exist we have kept it in here to make the specification as self-contained as possible.

³⁰ Available under <http://vo-urp.googlecode.com/svn/trunk/uml/IVOA%20UML%20Profile%20v-3.xml>

This is the root of the complete model, contains all packages, classes etc. Also contains any imported profile.

Stereotypes

- **<<model>>**

If the designer wants to annotate the model with the tags in this stereotype (s)he must explicitly associate this stereotype to the Model.

Tags:

- **author** : Indicates the author(s) of the model.
- **title** : provides a long title to the model. The name of the model is assumed to be short.
- **subject**: 0..* list of subjects in the sense of the Registry's subject attribute.

B.3 Package

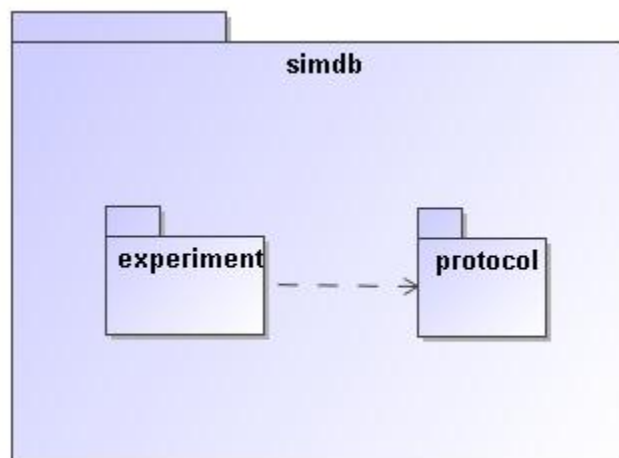


Figure 12: This figure shows a package "simdb" that contains two other packages. Of these the experiment package depends on the protocol packages, which is indicated by the dashed arrow. See Figure 2 for the somewhat more complex package structure used in SimDM.

A package groups related elements such as class definitions and possibly sub packages. Packages can depend on each other (indicated by the dashed line), which means that elements in one package can use elements in the target package in their definition. This relation is transitive. A package is similar to an XML namespace and in fact we map UML packages to XML namespaces in the XML schema mapping for the model described in 4.2.

B.4 Class



Figure 13: A Class is a rectangular box, with the name of the class in boldface.

Classes are the fundamental building blocks of a data model. A Class represents a full-fledged concept and is built up from properties and relations to other Classes. An important feature of Classes as opposed to DataTypes (see below) is that instances of Classes, i.e. objects, have their own, explicit identity³¹. That is we want to assign an explicit identifier to each concept so that

Properties:

- **isAbstract**
Indicated by *italicised* name of the object. Implies that no instances can be made of the class, only of concrete (=non abstract) sub classes.

B.5 *ValueType*

A ValueType represents a simple concept that is used to describe/define more complex concepts such as Classes. ValueType-s are, in contrast to Classes not separately identified. They are identified by their value. For example an integer is a value type; all instances of the integer value 3 represent the same integer. In this profile ValueType-s are only represented using specialised examples. Attributes (see below) must have a ValueType as their datatype.

B.6 *PrimitiveType*

PrimitiveTypes are the simplest examples of ValueTypes. They are represented by a single value only. A set of PrimitiveTypes is predefined in the IVOA profile (see Figure 14).

³¹ This is admittedly a somewhat theoretical object-oriented concept,

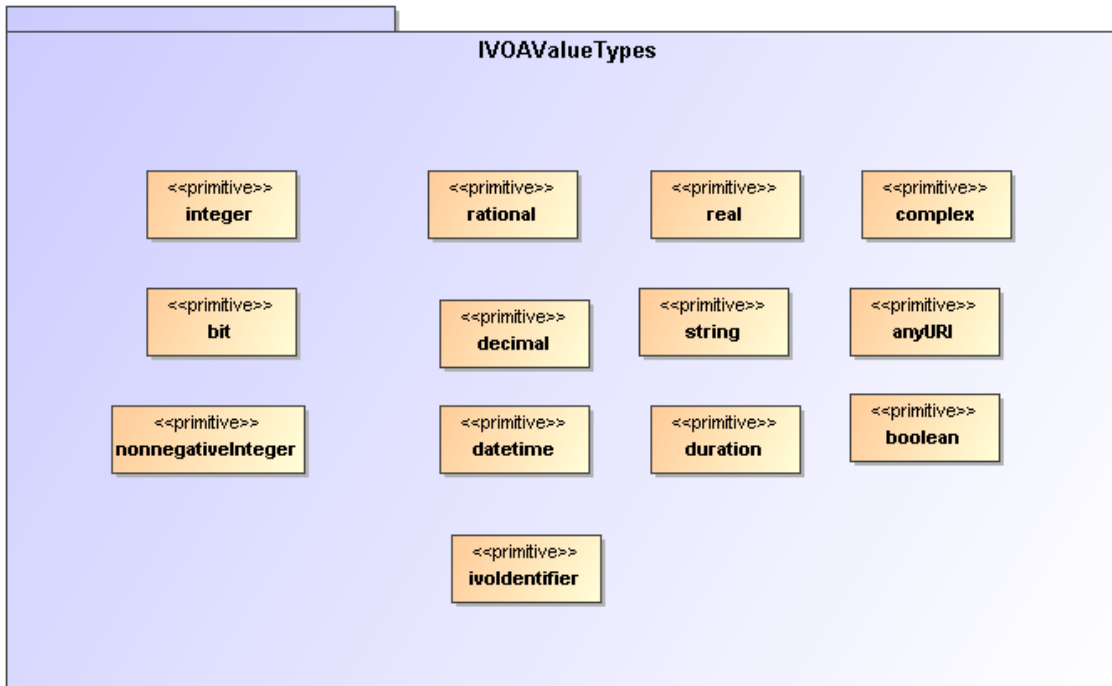


Figure 14: The PrimitiveTypes that are predefined in the IVOA profile.

B.7 *DataType*

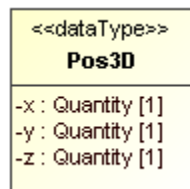


Figure 15: Example of a structured datatype: Pos3D represents a position in 3D space and is defined using x, y and z attributes. The DataType symbol is distinguished from the Class by the <<dataType>> stereotype.

A DataType is a ValueType that has more structure than a single value. This structure is modelled using Attributes, just as on ObjectTypes.

B.8 *Enumeration*

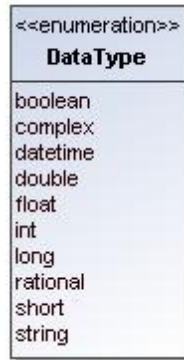


Figure 16: An enumeration is indicated by a box with the name of the the enumeration and the list of literals.

An Enumeration is a ValueType that is defined by a list of valid values. These are the only values that instances of this data type can assume.

B.9 Attribute

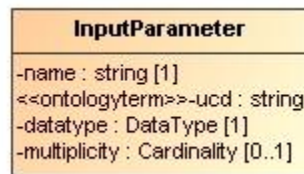


Figure 17: An attribute is indicated by a line with a name, a datatype, an indication of the multiplicity and possibly a stereotype.

An Attribute is a Property of a type (object type as well as structured data type). An attribute's data type is always a Value type, not an object type. For object type properties one should use [References](#)

Properties

- data type
- multiplicity/cardinality: indicates the cardinality of the attribute (assumed to be 0..1, or 1. This is a relational bias based on normal form and the assumption that most databases do not allow storage of arrays in single columns.)

Stereotypes

- <<attribute>>
 - To assign further properties such as the tags this stereotype attribute must be explicitly assigned.
 - Tag definitions
 - length [integer]: Constraint indicating that an attribute must have a specific fixed length. Is relevant only for attributes of type string.

- maxLength [integer]: Constraint indicating that an attribute may at most have the indicated length. Is relevant only for attributes of type string. Is used in mappings to TAP to indicate the length of the corresponding column. This would seem to be very much an application specific feature and therefore belong to logical modelling. But this profile can be used for that purpose, hence it is included.
- uniqueGlobally: Constraint indicating that only one instance of the type of the Class owning this attribute can have a given value. Globally should be read to mean globally in a given instance of the model, i.e. a database for example that stores instances of the model.
- uniqueInCollection [boolean]: If true, indicates that the value of the attribute can not be shared by the same attribute of any other instance of the Class owning this attribute that is in the same collection, i.e. has the same container object. In SimDB/DM an example is given by the name attribute of the InputParameter class. For a given Protocol
- <<ontologyterm>>

There are many instances in the data model where we need to describe elements of the SimDB/Resource-s explicitly, because we do not have implicit information based on the context. Examples are the various properties of object types, the target objects and processes etc. Apart from a name and a description we then frequently add an attribute which is supposed to "label" the element according to an assumed standard list of terms.

We model this using the <<ontologyterm>> stereotype. Attributes with this stereotype are assumed to take their values from such a predefined "ontology"³².

Tag definitions:

 - ontologyURI

A URL locating a standard (RDF|SKOS|OWL|???) document containing a list of terms from which the value for this attribute may be obtained. It is our opinion that the Semantics working group should be responsible for the definition of relevant ontologies (or semantic vocabularies, or thesauri, or ...) required for a given application domain, though the contents should be decided in cooperation with domain experts.

B.10 Inheritance

³² Possibly this should be a vocabulary, that at least is intended, and the stereotype might have to be called <<skosterm... with tag definition named skosVocabulary.

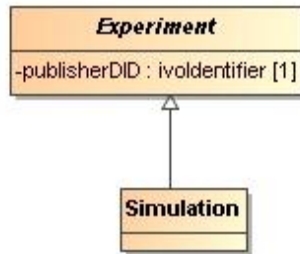


Figure 18: Inheritance is indicated by a line with an open arrow from a subclass to its base class.

Indicates the typical “is a” relation between the sub-class and its base-class (the one pointed at). In this profile we do not support multiple inheritance.

B.11 Collection

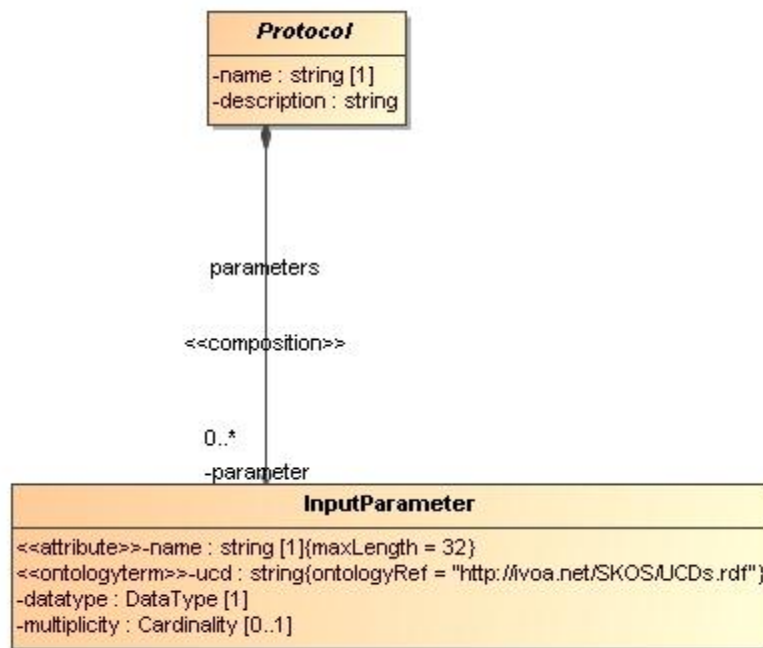


Figure 19: The line with the closed circle on one end and an arrow on the other indicates a composition relation, or collection, between the parent (on the side of the circle) and the child, on the other side.

This relation indicates a composition relation between one, parent object and 0 or more child objects. The life cycles of the child objects are governed by that of the parent.

In UML a composition relation is represented by a binary association end.

B.12 Reference

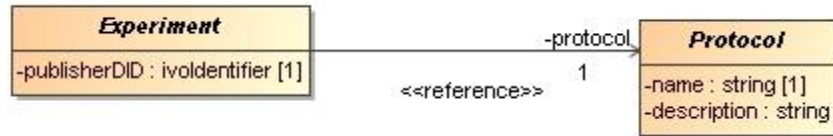


Figure 20: A Reference is represented by a line connecting a class with another, referenced, class with an arrow on the referenced class. Note, the <<reference>> stereotype indication is not required.

This is a relation that indicates a kind of usage, or dependency of one object on another. It is in general shared, i.e. many objects may reference a single other object. Accordingly the referenced object is independent of the "referee". In our profile the cardinality can not be > 1.

For implementing the Reference in UML we use a shared, navigable binary association end.

B.13 Subsets

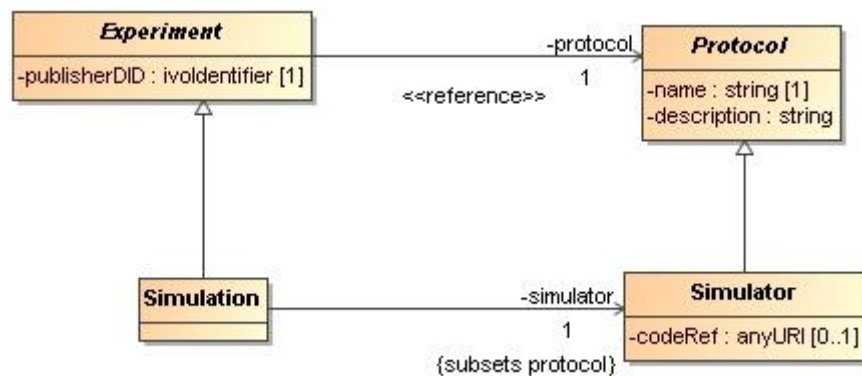


Figure 21: The subsets property can only be assigned to a relation between two objects that are both subclasses of Classes that have an equivalent relation. It is indicated by the {subsets ...} annotation to the relevant association end.

The "subsets" property, when associated to a Reference or Collection (in UML to the corresponding association end), indicates that a relation overrides the definition of a relation of the same name defined on a base class. It does so by specifying that the class at the end point of the relation should be a subclass of the class at the endpoint of the original, sub-setted relation.

Appendix C Issues

We identify and discuss here several issues with the SimDM that need to be discussed further.

C.1 Normalisation

The current version of the SimDB/DM is rather more *normalised* [9] than most of the other data models in the IVOA. We explain this concept based on a particular choice we made during the modelling process, and then we discuss the consequences of particular choices.

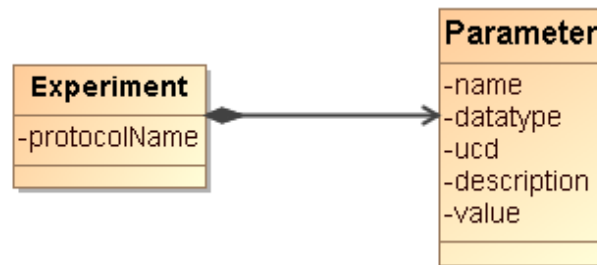


Figure 22: Non-normalised model for experiments and parameters.

At an earlier stage, the model was less normalised in the design of the input parameters of an experiment, as is illustrated in Figure 22. There was no separate protocol class, only an attribute *protocolName* on the **Experiment** class indicated the protocol by which the experiment was run. Also, the input parameters on the experiment were completely contained in a collection of **Parameter**-s. The **Parameter** class contained all the details, including *name* of the parameter, *description*, *ucd* etc. It also contained the *value* of the parameter in the experiment.

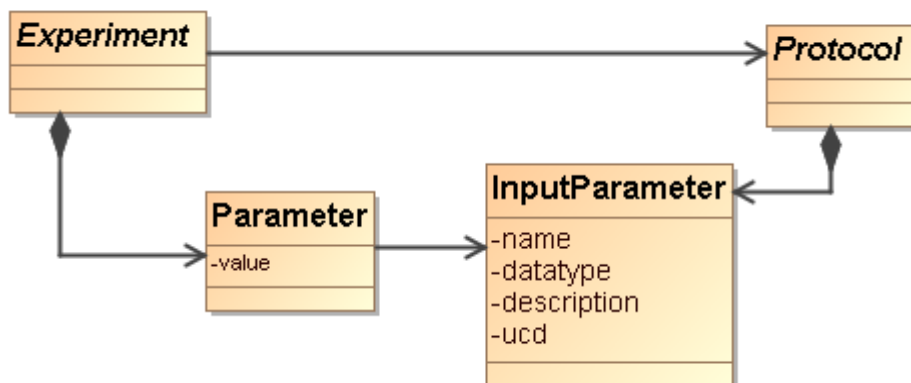


Figure 23: Normalised modelling of experiments, protocols and their parameters.

Currently the model treats parameter definitions and settings as in Figure 23. In this normalised design, the *Protocol* is given a class of its own, and it contains the input parameter collection. The *InputParameter* class does *not* contain a

value, only the definition of the parameter: *name*, *datatype*, *ucd*, *description*. The values assigned to parameters in a given experiment are captured with the *ParameterSetting* class, contained in a collection off *Experiment*.

The motivation for this change of model was that a SimDB instance will in general contain many simulations (experiments) run with the same simulator code (protocol). In the old model, each experiment has to define the collection of input parameters with all details. In the new design this only has to be defined once, on the appropriate protocol. This clearly is less redundant, which is one important design goal of a normalised modelling approach. At the same time it provides an explicit identity to input parameters, which allows us to ask explicit questions about all parameter settings for a given, identified parameter. In the old model this is only indirectly possible, using equality of the name of input parameters for all experiments having the same *protocolName*. Now we can ask for all experiments with the same protocol reference, and look for parameter settings with the same input parameter reference. This is arguably a more "correct" model of reality.

There are therefore advantages to normalisation, but there are also disadvantages. We need to realise these and make choices that optimise the usability of the data model. One of the main disadvantages is that *references*, which naturally have to be introduced when normalising a model, are more difficult to deal with than most of the other modelling elements, particularly in some physical representations (see below). When defining a new experiment, one will have to find the input parameter that one needs to set, and instead of simply giving name/value, one needs to represent the reference to the parameter. For this one may have to extract the protocol as stored in the SimDB and find the appropriate identifiers of the input parameters. In this sense an *Experiment* definition becomes less self-contained. It depends on the details of the *registered Protocol*. This protocol is registered separately and necessarily at an earlier point in time.

This puts strong requirements on SimDB implementations to maintain referential integrity, something which will be even harder to achieve if we were to allow cross-SimDB referencing. In one advanced usage scenario the UC San Diego version of SimDB registers the Enzo³³ simulator, whilst the Italian SimDB allows registration of simulations that used it and reference the remote protocol³⁴.

Similarly a query language needs to be able to handle with this level of indirection. For example in a relational database one needs to write joins between *ParameterSetting* and *InputParameter*. For expert SQL users this is not a problem, but is something to get used to. For simpler query languages, those not allowing joins, like TAP/Param, asking meaningful queries becomes very difficult. One way around this problem could be to add some view definitions to the model. In relational databases, views are predefined, named SQL queries that can be treated as if they were tables when querying the database. It is quite straightforward to define some SQL queries that as it were denormalise the

³³ <http://lca.ucsd.edu/portal/software/enzo>

³⁴ We actually do not support this scenario in the current version of SimDB.

model and put the input parameter definition back under the experiment together with the value. This way one may protect users of the database from the high level of normalisation.

C.2 Quantities and Units

At various locations in the SimDM data model numerical values can be defined, for example in parameter settings or the characterisation of properties of representation object type collections. Often these numerical values will need to have a unit. The IVOA has two ways of dealing with units. Either units are fixed explicitly for properties/parameters in protocol or data model, sometimes depending on the small list of possible UCDs. Alternatively units are explicitly stated, for example in VOTable. At the moment we support the second mode, especially because, as is true for VOTable, we do not know what kind of property is being used. To this end we introduce a value type in the model, `Quantity`, which contains a value and a unit, and which is the data type of various value attributes, for example in `NumericalParameterSetting` or `Characterisation`. In the XML schema this is translated to a complexType with 2 elements, in the relational database schema to two columns, one with the value, one with the unit. It is in the use of the relational schema that we anticipate problems with this approach, especially in the query protocol to SimDB. Consider the typical science question: return all N-body simulations with particle mass roughly $10^{10}M_{\text{sun}}$. In SimDB this would need to be translated in an ADQL query which contains the unit column explicitly. Allowing users' freedom of registering SimDB resources using any units they desire can lead to resource containing, for the same observable "N-body particle mass", values with a whole range of units. To provide reasonable support to users requires the SimDB implementation to be able to do the automated transformation. But in the We propose in SimDB to use ADQL/TAP as the query interface. If units are stored explicitly users can phrase queries using these,

An alternative approach is to mandate stating values for properties with a given UCD (or other semantic label) always with the same units. This would solve the query problem but poses others. For one it may be very (too?) unnatural for users to be forced to use meters for cosmological simulations, or megaparsec for simulations in the solar system. Related to this is the probably contentious discussion of what units to assign to what UCD. One might choose SI or cgs units, but these are not always very useful or natural.

C.3 Linking services, experiments and other resources.

When studying the SimDM and comparing it to the explanation in 2.4, one will notice that there is no concept of *storage* for the results in the proposed model. The reason is that whereas we can define the concept of a Result as collections of Objects quite satisfactory, we have shied away of trying to model the precise way these results are actually stored in files or a database. There are simply too

many possible and actual ways in which results can be stored in a file system. In general, a result, or snapshot in our model, cannot be modelled with a simple reference to a file or table in which it is stored. Large results may be split up over many files, stored as structs, or in arrays. We have therefore decided not to open this can of worms (or reopen it, see the Quantity data model <http://www.ivoa.net/internal/IVOA/IvoaDataModel/qty23.pdf>).

Instead we assume the existence of web services that allow users access to the results of SimDB experiments. Some of these services may implement a standard protocol as defined by SimDAP, or they may be custom services. The precise way to relate experiments to services and what can be inferred about how to call them is the task of the SimDAP protocol.

In the model actually web services are related to resources in general. This can be used to represent services that gives access to a set of experimental results from some project, or that can for example visualise any result of experiments performed according to a fixed protocol, for example generic Gadget-format visualisers.

Appendix D Use of SimDB data model

Here we present our views on how the data model presented in this and accompanying documents could/should be used. We first discuss the usage in the SimDB service protocol for which it was originally intended. We then discuss the possible usage in SimDAP that was assumed to use the model in some form as well, though the link has been weakened after the SNAP proposal was separated in SimDB and SimDAP. We then propose how also the S3 proposal [14] could use the model, thus tying that specification closer to the other developments in the theory interest group. Finally we discuss how one might identify Registry Resource-s in the model, and how one might extract these, in effect creating a specialised Registry implementation on top of the Simulation Database.

D.1 *SimDB*

As will be discussed in the *SimDB Services Note* [in preparation], the current proposal is that a Simulation Database uses the SimDB data model exactly as described here, in all its details. With one proviso though. The data model presented here is in UML form, which is not usable in protocols, databases etc. Hence one or more physical representations must be created. The XML Schema representation was shortly discussed in section 4.2 and will be more detailed in the SimDB note. But the main emphasis in the SimDB protocol is on a TAP/ADQL interface, which requires a relational database and hence a relational model. In the SimDB note we therefore provide a relational mapping of the data model together with its representation in metadata formats prescribed by the TAP specification.

This means that, in the current proposal, the data model only needs a transformation to a physical representation, without any loss of information. However as we have discussed in section C.1, the highly normalised nature of the model complicates its usage, whatever the representation. There we proposed a solution using “denormalised representations” of the model. These representations could be created in UML and used exactly as the current model as the basis for UTYPEs, XML schema, TAP metadata etc.

We believe that in principle it is up to applications³⁵ how to *use* the data model, but we feel that it would be good if some more or less formal transformation were given that could express the application model into the SimDB model, vice versa transform the SimDB model into the application specific representation.

D.2 *SimDAL/SimTAP*

The Simulation Data Access Layer (SimDAL) “... defines a protocol for retrieving data coming from numerical simulations.” When SNAP was split into two

³⁵ This includes Protocols such as SimDAP and S3, but also SimDB itself. And it includes anyone who wishes to give an custom web (service) interface to their SimDB installation.

separate efforts, SimDAP was supposed to be the data access component, complementing the discovery functionality from SimDB. The main tasks of SimDAP were how to handle accessing the possibly very large data sets resulting from 3+1D simulations. In the subsequent generalisation to other types of simulations and models, SimDAP was transformed to SimDAL.

This should be a “regular” DAL-type protocol, with a queryData component. That component is supposed to be basically a TAP service on a data model derived from SimDM. The details how this transformation is accomplished are described in the SimDAL document [22]