

# CUAHSI Community Observations Data Model (ODM) Version 1.1 Design Specifications

May 2008

David G. Tarboton<sup>1</sup>, Jeffery S. Horsburgh<sup>1</sup>, David R. Maidment<sup>2</sup>

## Abstract

The CUAHSI Hydrologic Information System project is developing information technology infrastructure to support hydrologic science. One aspect of this is a data model for the storage and retrieval of hydrologic observations in a relational database. The purpose for such a database is to store hydrologic observations data in a system designed to optimize data retrieval for integrated analysis of information collected by multiple investigators. It is intended to provide a standard format to aid in the effective sharing of information between investigators and to allow analysis of information from disparate sources both within a single study area or hydrologic observatory and across hydrologic observatories and regions. The observations data model is designed to store hydrologic observations and sufficient ancillary information (metadata) about the data values to provide traceable heritage from raw measurements to usable information allowing them to be unambiguously interpreted and used. A relational database format is used to provide querying capability to allow data retrieval supporting diverse analyses. A generic template for the observations database is presented. This is referred to as the Observations Data Model (ODM).

## Introduction

The Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI) is an organization representing more than 100 universities and is sponsored by the National Science Foundation to provide infrastructure and services to advance the development of hydrologic science and education in the United States. The CUAHSI Hydrologic Information System (HIS) is being developed as a geographically distributed network of hydrologic data sources and functions that are integrated using web services so that they function as a connected whole. One aspect of the CUAHSI HIS is the development of a standard database schema for use in the storage of point observations in a relational database. This is referred to as the point Observations Data Model (ODM) and is intended to allow for comprehensive analysis of information collected by multiple investigators for varying purposes. It is intended to expand the ability for data analysis by providing a standard format to share data among investigators and to facilitate analysis of information from disparate sources both within a single study area or hydrologic observatory and across hydrologic observatories and regions. The ODM is designed to store hydrologic observations with sufficient ancillary information (metadata) about the data values to provide traceable heritage from raw measurements to usable information allowing them to be unambiguously interpreted and used. Although designed specifically with hydrologic observation data in mind, this data model has a simple and general structure that will also

---

<sup>1</sup> Utah Water Research Laboratory, Utah State University

<sup>2</sup> Center for Research in Water Resources, University of Texas at Austin

accommodate a wide range of other data, such as from other environmental observatories or observing networks.

ODM uses a relational database format to allow for ease in querying and data retrieval in support of a diverse range of analyses. Reliance on databases and tables within databases also provides the capability to have the model scalable from the observations of a single investigator in a single project through the multiple investigator communities associated with a hydrologic observatory and ultimately to the entire set of observations available to the CUAHSI community. ODM is focused on observations made at a point. A relational database model with individual observations recorded as individual records (an atomic model) has been chosen to provide maximum flexibility in data analysis through the ability to query and select individual observation records. This approach carries the burden of record level metadata, so it is not appropriate for all variables that might be observed. For example, individual pixel values in large remotely sensed images or grids are inappropriate for this model.

This data model is presented as a generic template for a point observations database, without reference to the specific implementation in a database management system. This is done so that the general design is not limited to any specific proprietary software, although we expect that implementations will take advantage of capabilities of specific software. It should be possible to implement ODM in a variety of relational database management systems, or even in a set of text tables or variable arrays in a computer program. However, to take full advantage of the relationships between data elements, the querying capability of a relational database system is required. By presenting the design at a general conceptual level, we also avoid implementation specific detail on the format of how information is represented. See the discussion of Dates and Times under ODM features below for an example of the distinction between general concepts and implementation specific details.

### **Version Information**

ODM has evolved from an initial design presented at a CUAHSI workshop held in Austin during March, 2005 (Maidment, 2005) that was then widely reviewed with comments being received from 22 individuals (Tarboton, 2005). These reviews served as the basis for a redesign that was presented at a CUAHSI workshop in Duke during July, 2005 and presented as part of the CUAHSI HIS status report (Horsburgh et al., 2005). Following this presentation of the design, the data model was reviewed and commented on by a number of others, including the CLEANER (Collaborative Large-scale Engineering Analysis Network for Environmental Research) cyberinfrastructure committee. Further versions of the Observations Data Model were circulated in April, June and October 2006. These documented changes made in the evolution of this design. The fundamental design, however, has not changed since the status report presentation of the model (Horsburgh et al., 2005) but many table and field names have been changed. Tables have also been added to give spatial reference information, metadata information, and to define controlled vocabularies. Version 1.0 of ODM, which was the first release version of ODM, has been implemented and tested within the WATERS network of test bed sites and was documented in Water Resources Research (Horsburgh et al., 2008). This document describes the second release version of the data model design, which has been named ODM Release Version 1.1, and has been so named to correspond to the Version 1.1 release of the CUAHSI HIS. This document supersedes the previous documents.

In general, the following changes have been made for Version 1.1:

- All integer IDs serving as the primary key for tables in ODM have been changed to auto number/identity fields.
- Text field lengths have been relaxed in some cases and have been standardized according to the following scheme: codes = 50 characters, terms = 255 characters, links = 500 characters, definitions/explanations = unlimited.
- Check constraints have been defined for the Latitude and Longitude fields in the Sites table.
- Check constraints have been added to many of the fields in ODM to constrain the characters that are valid for those fields (see Appendix A for details).
- Relationships have been added between controlled vocabulary tables and the tables that contain the fields that they define. This was done to more rigorously enforce the ODM controlled vocabularies.
- Unique constraints were placed on both SiteCode in the Sites table and VariableCode in the Variables table.
- The controlled vocabulary was relaxed on the QualityControlLevels table to allow more detailed versioning of data series. A QualityControlLevelCode was also added to this table to facilitate this.
- A Citation field was added to the Sources table to provide a place for a formal citation for data in the database.
- A Speciation field was added to the Variables table. This provides a place to store information about the speciation of chemistry observations. A SpeciationCV controlled vocabulary table was added to define this field.
- An ODMVersion table was added to store the version number of the database.
- The SeriesCatalog table has been updated based on the addition of the above fields.

## Hydrologic Observations

Many organizations and individuals measure hydrologic variables such as streamflow, water quality, groundwater levels, and precipitation. National databases such as USGS' National Water Information System (NWIS) and USEPA's data Storage and Retrieval (STORET) system contain a wealth of data, but, in general, these national data repositories have different data formats, storage, and retrieval systems, and combining data from disparate sources can be difficult. The problem is compounded when multiple investigators are involved (as would be the case at proposed CUAHSI Hydrologic Observatories) because everyone has their own way of storing and manipulating observational data. There is a need within the hydrologic community for an observations database structure that presents observations from many different sources and of many different types in a consistent format.

Hydrologic observations are identified by the following fundamental characteristics:

- The location at which the observations were made (space)
- The date and time at which the observations were made (time)
- The type of variable that was observed, such as streamflow, water surface elevation, water quality concentration, etc. (variable)

These three fundamental characteristics may be represented as a data cube (Figure 1), where a particular observed data value (D) is located as a function of where it was observed (L), its time of observation (T), and what kind of variable it is (V), thus forming  $D(L,T,V)$ .

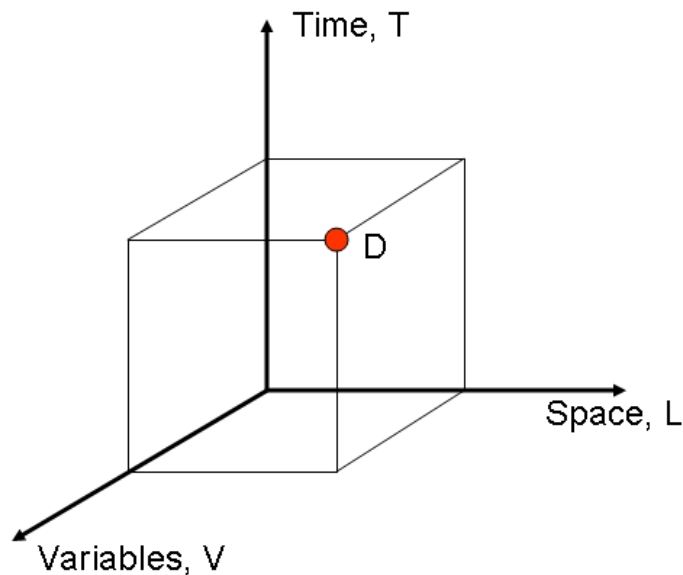


Figure 1. A measured data value (D) is indexed by its spatial location (L), its time of measurement (T), and what kind of variable it is (V).

In addition to these fundamental characteristics, there are many other distinguishing attributes that accompany observational data. Many of these secondary attributes provide more information about the three fundamental characteristics mentioned above. For example, the

location of an observation can be expressed as a text string (i.e., “Bear River Near Logan, UT”), or as latitude and longitude coordinates that accurately delineate the location of the observation. Other attributes can provide important context in interpreting the observational data. These include data qualifying comments and information about the organization that collected the data. The fundamental design decisions associated with the ODM involve choices as to how much supporting information to include in the database and whether to store (and potentially repeat) this information with each observation or save this information in separate tables with key fields used to logically associate observation records with the associated information in the ancillary tables. Table 1 presents the general attributes associated with a point observation that we judged should be included in the generic ODM design.

**Table 1. ODM attributes associated with an observation**

<b>Attribute</b>	<b>Definition</b>
Data Value	The observation value itself
Accuracy	Quantification of the measurement accuracy associated with the observation value
Date and Time	The date and time of the observation (including time zone offset relative to UTC and daylight savings time factor)
Variable Name	The name of the physical, chemical, or biological quantity that the data value represents (e.g. streamflow, precipitation, temperature)
Speciation	For concentration measurements, the species in which the concentration is expressed (e.g., as N, or as NO <sub>3</sub> , or as NH <sub>4</sub> )
Location	The location at which the observation was made (e.g. latitude and longitude)
Units	The units (e.g. m or m <sup>3</sup> /s) and unit type (e.g. length or volume/time) associated with the variable
Interval	The interval over which each observation was collected or implicitly averaged by the measurement method and whether the observations are regularly recorded on that interval
Offset	Distance from a reference point to the location at which the observation was made (e.g. 5 meters below water surface)
Offset Type/ Reference Point	The reference point from which the offset to the measurement location was measured (e.g. water surface, stream bank, snow surface)
Data Type	An indication of the kind of quantity being measured (e.g. a continuous, minimum, maximum, or cumulative measurement)
Organization	The organization or entity providing the measurement
Censoring	An indication of whether the observation is censored or not
Data Qualifying Comments	Comments accompanying the data that can affect the way the data is used or interpreted (e.g. holding time exceeded, sample contaminated, provisional data subject to change, etc.)
Analysis Procedure/ Method	An indication of what method was used to collect the observation (e.g. dissolved oxygen by field probe or dissolved oxygen by Winkler Titration) including quality control and assurance that it has been subject to
Source	Information on the original source of the observation (e.g. from a specific organization, agency, or investigator 3 <sup>rd</sup> party database)
Sample Medium	The medium in which the sample was collected (e.g. water, air, sediment, etc.)
Value Category	An indication of whether the data value represents an actual measurement, a calculated value, or is the result of a model simulation

## Observations Data Model

The schema of the Observations Data Model is given in Figure 2. Appendix A gives details of each table and each field in this generic data model schema. Appendix A serves as the data dictionary for the data model and documents specific database constraints, data types, examples, and best practices. The primary table that stores point observation values is the DataValues table at the center of the schema in Figure 2. Logical relationships between fields in the data model are shown and serve to establish the connectivity between the observation values and associated ancillary information. Details of the relationships are given in Table 2. Figure 2 shows each of the controlled vocabulary tables and their relationships to the table containing the field that they define. Controlled vocabulary tables are highlighted with red headers. In Figure 2, each of the mandatory fields is shown in bold text, whereas optional fields are shown in regular text.

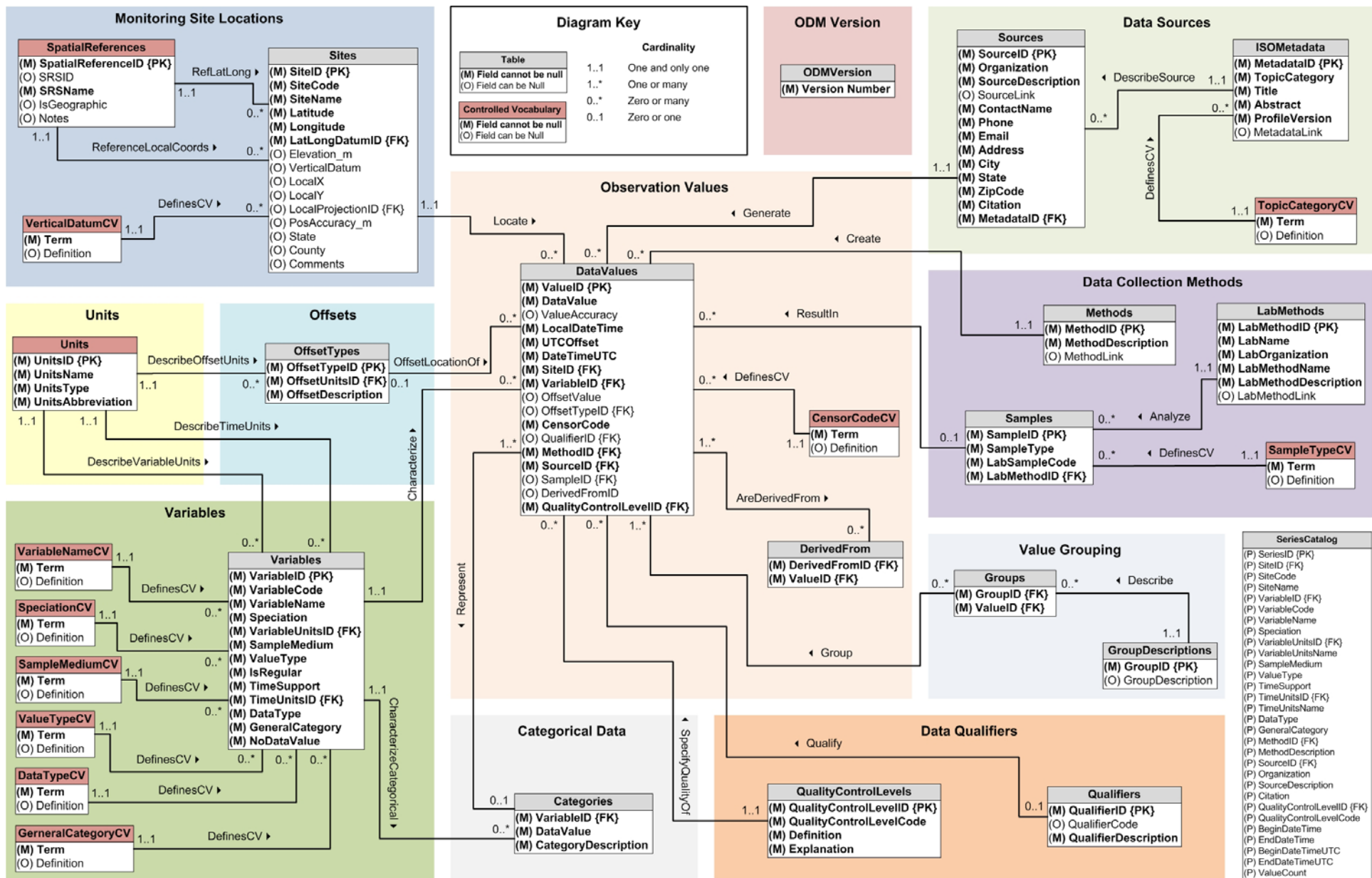


Figure 2. Observations Data Model schema.

Table 2. Observations Data Model Logical Relationships

Relationships that define <u>ancillary information</u> about data values				
Table	Field	Type	Field	Table
DataValues	SiteID	* <-> 1	SiteID	Sites
DataValues	VariableID	* <-> 1	VariableID	Variables
DataValues	OffsetTypeID	* <-> 1	OffsetTypeID	OffsetTypes
DataValues	QualifierID	* <-> 1	QualifierID	Qualifiers
DataValues	MethodID	* <-> 1	MethodID	Methods
DataValues	SourceID	* <-> 1	SourceID	Sources
DataValues	SampleID	* <-> 1	SampleID	Samples
DataValues	QualityControlLevelID	* <-> 1	QualityControlLevelID	QualityControlLevels

Relationships that define <u>derived from groups</u>				
Table	Field	Type	Field	Table
DataValues	DerivedFromID	* <-> *	DerivedFromID	DerivedFrom
DataValues	ValueID	1 <-> *	ValueID	DerivedFrom

Relationships that define <u>groups</u>				
Table	Field	Type	Field	Table
DataValues	ValueID	1 <-> *	ValueID	Groups
GroupDescriptions	GroupID	1 <-> *	GroupID	Groups

Relationships used to define <u>categories for categorical data</u>				
Table	Field	Type	Field	Table
Variables	VariableID	1 <-> *	VariableID	Categories
DataValues	DataValue	* <-> 1	DataValue	Categories

Relationships used to define the <u>Units</u>				
Table	Field	Type	Field	Table
Units	UnitsID	1<->*	VariableUnitsID	Variables
Units	UnitsID	1<->*	TimeUnitsID	Variables
Units	UnitsID	1<->*	OffsetUnitsID	OffsetTypes

Relationship used to define the <u>Sample Laboratory Methods</u>				
Table	Field	Type	Field	Table
LabMethods	LabMethodID	1<->*	LabMethodID	Samples

Relationships used to define the <u>Spatial References</u>				
Table	Field	Type	Field	Table
SpatialReferences	SpatialReferenceID	1<->*	LatLongDatumID	Sites
SpatialReferences	SpatialReferenceID	1<->*	LocalProjectionID	Sites

Relationship used to define the <u>ISOMetaData</u>				
Table	Field	Type	Field	Table
IsoMetaData	MetadataID	1<->*	Sources	MetadataID



Relationships used to define <u>Controlled Vocabularies</u>				
Table	Field	Type	Field	Table
VerticalDatumCV	Term	1<->*	VerticalDatum	Sites
SampleTypeCV	Term	1<->*	SampleType	Samples
VariableNameCV	Term	1<->*	VariableName	Variables
ValueTypeCV	Term	1<->*	ValueType	Variables
DataTypeCV	Term	1<->*	DataType	Variables
SampleMediumCV	Term	1<->*	SampleMedium	Variables
SpeciationCV	Term	1<->*	Speciation	Variables
GeneralCategoryCV	Term	1<->*	GeneralCategory	Variables
TopicCategoryCV	Term	1<->*	TopicCategory	ISOMetadata
CensorCodeCV	Term	1<->*	CensorCode	DataValues

Relationship type is indicated as One to One (1<->1), One to Many (1<->\*), Many to One (\*<->1) and Many to Many (\*<->\*). The first set of relationships defines the links to tables that contain ancillary information. They are used so that only compact (integer) identifiers are stored with each data value and thus repeated many times while the more voluminous ancillary information is stored to the side and not repeated. The second set of relationships defines derived from groupings used to specify data values that have been used to derive other data values. The third set of relationships defines logical groupings of data values. The fourth set of relationships is used to specify the categories associated with categorical variables. The fifth set of relationships is used to define the units. The sixth set of relationships associates laboratory methods with samples. The seventh set of relationships associates sites with the Spatial Reference System used to define the location. The eighth set of relationships associates project and dataset level metadata with each data source. The last set of relationships defines the linkage between the controlled vocabulary fields and the tables that stored the acceptable terms for those fields. Details of how these relationships work are given in the discussion of features of the data model design below.

## Features of the Observations Data Model Design

### Geography

ODM is intended to be independent of the geographical representation of the site locations. The geographic location of sites is specified through the Latitude, Longitude, and Elevation information in the Sites table, and optionally local coordinates, which may be in a standard geographic projection for the study area or in a locally defined coordinate system specific to a study area. Each site also has a unique identifier, SiteID, which can be logically linked to one or more objects in a Geographic Information System (GIS) data model. For example, Figure 3 depicts a one-to-one relationship between sites within ODM and HydroPoints within the Arc Hydro Framework Data Model (Maidment, 2002) used to represent objects in a digital watershed. In simple implementations, SiteID may have the same integer value as the identifier for the associated GIS object, HydroID in this case. In more complex implementations, and especially when multiple databases are merged into a single ODM, it may not be possible to preserve the simple one-to-one relationship between SiteID and HydroID with each of these fields holding the same integer identifier values. In these cases, where SiteID and HydroID are

not the same, a coupling table would be used to associate the ODM SiteIDs used to identify sites with HydroIDs in the Arc Hydro data model.

SiteID must be unique within an instance of ODM. This could, for example, be achieved by assigning SiteIDs from a master table. The linkage between SiteIDs and GIS object IDs is intended to be generic and suitable for use with any geographic data model that includes information specifying the location of sites. For example, a linear referencing system on a river network, such as the National Hydrography Dataset, might be used to specify the location of a site on a river network. Addressing relative to specific hydrologic objects through the SiteID field provides direct and specific location information necessary for proper interpretation of data values. Information from direct addressing relative to hydrologic objects is often of greater value to a user than the simple Latitude and Longitude information stored in the ODM Sites table. For example, it is more useful to know that a stream gage is on such and such a stream rather than simply its latitude and longitude.

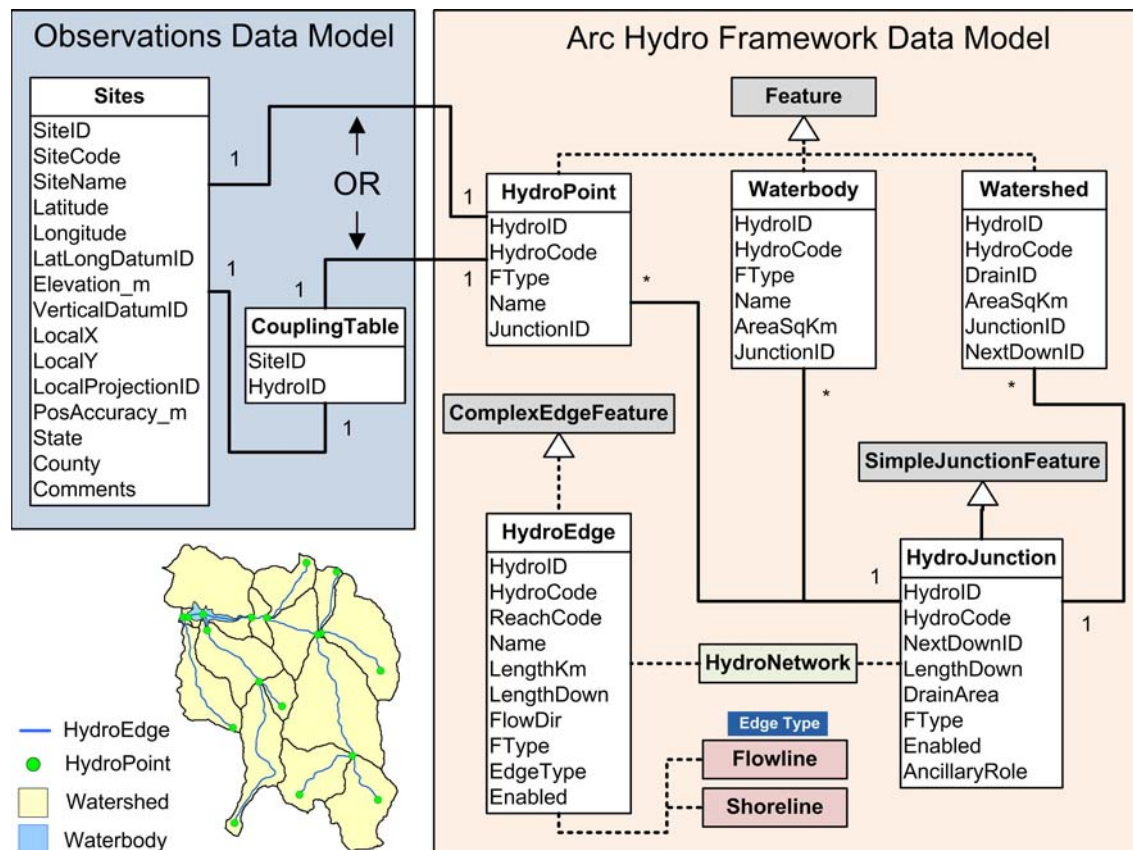


Figure 3. Arc Hydro Framework Data Model and Observations Data Model related through SiteID field in the Sites table.

### Series Catalog

A "data series" is an organizing principle used in ODM. A data series consists of all the data values associated with a unique site, variable, method, source, and quality control level combination in the DataValues table. The SeriesCatalog table lists data series, identifying each by a unique series identifier, SeriesID. This table is essentially a summary of many of the tables

in the ODM and is not required to maintain the integrity of the data. However, it serves to provide a listing of all the distinct series of data values of a specific variable at a specific site. By doing so, this table provides a means by which users can execute most common data discovery queries (i.e., which variables have data at a site, etc.) without the overhead of querying the entire DataValues table, which can become quite large.

The SeriesCatalog table is also intended to support CUAHSI Web Service method queries such as GetSiteInfo, which returns information about a monitoring site within an instance of the ODM including the variables that have been measured at that site. It should be noted that data series, as they are defined here, do not distinguish between different series of the same variable at the same site but measured with different offsets. If for example temperature was measured at two different offsets by two different sensors at one site, both sets of data would fall into one data series for the purposes of the SeriesCatalog table. In these cases, interpretation or analysis software will need to specifically examine and parse the offsets by examining the offset associated with each data value. The SeriesCatalog table does not do this because its principal purpose is data discovery, which we did not want to be overly complicated. The SeriesCatalog table should be programmatically generated and modified as data are added to the database.

### Accuracy

Each data value in the DataValues table has an associated attribute called ValueAccuracy. This is a numeric value that quantifies the total measurement accuracy defined as the nearness of a measurement to the true or standard value. Since the true value is not known, the ValueAccuracy is estimated based on knowledge of the instrument accuracy, measurement method, and operational environment. The ValueAccuracy, which is also called the uncertainty of the measurement, compounds the estimates of both bias and precision errors. Bias errors are generally fixed or systematic and cannot be determined statistically, while precision errors are random, being generated by the variability in the measurement system and operational environment. Figure 4 illustrates the effects of these errors on a sample of measurements. Bias errors are usually estimated through specially designed experiments (calibrations). The precision errors are determined using statistical analysis by quantifying the measurement scatter, which is proportional to the standard deviation of the sample of repeated measurements. The total error is obtained by the root-sum-square of the estimates for bias and precision errors involved in the measurement. Figure 5 gives another illustration of the ValueAccuracy concept based on the analogy of a target, where the bulls eye at the center represents the true value.

ValueAccuracy is a data value level attribute because it can change with each measurement, dependent on the instrument or measurement protocol. For example, if streamflow is measured using a V-notch weir, it is actually the stage that is measured, with accuracy limited by the precision and bias of the depth recording instrument. The conversion to discharge through the stage-discharge relationship results in greater absolute error for larger discharges. Inclusion of the ValueAccuracy attribute, which will be blank for many historic datasets because historically accuracy has not been recorded, adds to the size of data in the ODM, but provides a way for factoring the accuracy associated with measurements into data analysis and interpretation, a practice that should be encouraged.

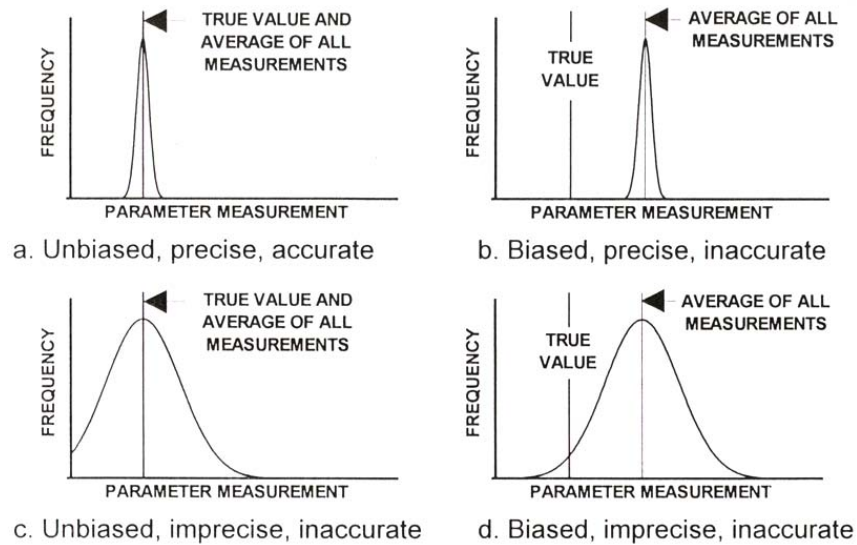


Figure 4. Illustration of measurement error effect (Source: AIAA, 1995).

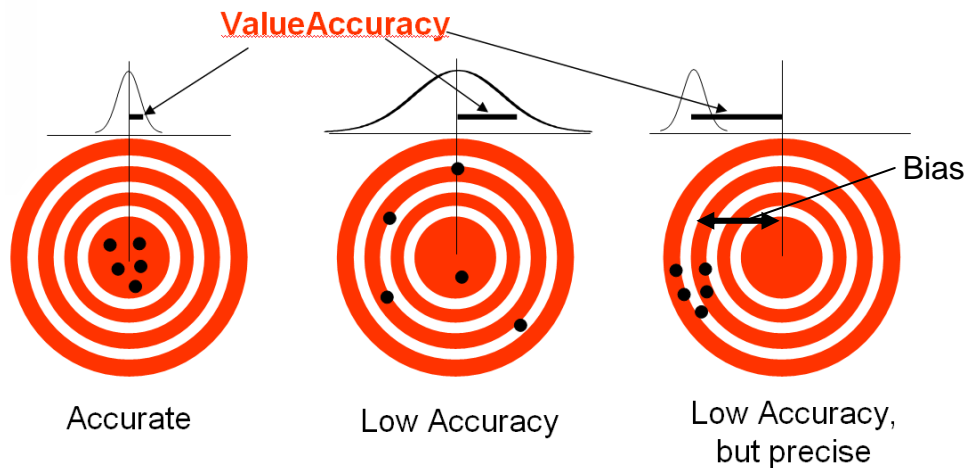


Figure 5. Illustration of Accuracy versus Precision (adapted from Wikipedia <http://en.wikipedia.org/wiki/Accuracy>).

In designing ODM, consideration was given to the suggestion by some reviewers to record bias and precision separately, in addition to ValueAccuracy for each data value. This has not been done at this release in the interest of parsimony and also because quantifying these separate components of the error is difficult. We suggest that for most measurements there should be the presumption that they are unbiased and that ValueAccuracy quantifies the precision and accuracy in the judgment of the investigator responsible for collecting the data. For cases where there is specific bias and precision information to complement the ValueAccuracy attribute, this could be recorded in the ODM as a separate variable, e.g. discharge precision, or temperature bias. The groups and derived from features (see below) could be used to associate these variables with their related observations. For measurements that are known to be biased, we suggest that the bias could be quantified by other reference measurements that should also be placed in the

database and that a new set of corrected measurements that have had the bias removed should be added to the database at a higher quality control level. These new measurements should have a lower ValueAccuracy value to reflect the improvement in accuracy by removal of the bias. The method and derived from information for these corrected measurements should give the bias removal method and refer to the data used to quantify and remove the bias.

### Offset

Each record in the DataValues table has two optional fields OffsetValue and OffsetTypeID. These are used to record the location of an observation relative to an appropriate datum, such as “depth below the water surface” or “depth below or above the ground.” The OffsetTypeID references an OffsetValue into an OffsetTypes table that gives units and definition associated with the OffsetValue. This design only has the capability to represent one offset for each data value. In cases (which we expect to be rare) when there are multiple offsets (e.g. distance in from a stream bank and depth below the surface) one of the offsets will need to be distinguished as a separate variable.

### Spatial Reference and Positional Accuracy

Unambiguous specification of the location of an observation site requires that the horizontal and vertical datum used for latitude, longitude, and elevation be specified. The SpatialReferences table is provided for this purpose to record the name and EPSG code of each Spatial Reference System used. EPSG codes are numeric codes associated with coordinate system definitions published by the OGP Surveying and Positioning Committee (<http://www.epsg.org/>). A non-standard Spatial Reference System, such as, for example, a local grid at an experimental watershed, may be defined in the SpatialReferences table Notes field. The accuracy with which the location of a monitoring site is known is quantified using the PosAccuracy\_m field in the Sites table. This is a numeric value intended to specify the uncertainty (as a standard deviation or root mean square error) in the spatial location information (latitude and longitude or local coordinates) in meters. Using a large number for PosAccuracy\_m (e.g. 2000 m) accommodates entry of data collected for a study area where the precise location where the observation was recorded is not known.

### Groups and Derived from Associations

The DerivedFrom and Groups tables fulfill the function of grouping data values for different purposes. These are tables where the same identifier (DerivedFromID or GroupID) can appear multiple times in the table associated with different ValueIDs, thereby defining an associated group of records. In the DerivedFrom table this is the sole purpose of the table, and each group so defined is associated with a record in the DataValues table (through the DerivedFromID field in that table). This record would have been derived from the data values identified by the group. The method of derivation would be given through the methods table associated with the data value. This construct is useful, for example, to identify the 96 15-minute unit streamflow values that go into the estimate of the mean daily streamflow. Note that there is no limit to how many groups a data value may be associated with, and data values that are derived from other data values may themselves belong to groups used to derive other data values (e.g. the daily minimum flow over a month derived from daily values derived from 15 minute unit values). Note also that a derived from group may have as few as one data value for the case where a data value is derived from a single more primitive data value (e.g., discharge from stage). Through this

construct the ODM has the capability to store raw observation values and information derived from raw observations, while preserving the connection of each data value to its more primitive raw measurement.

The GroupID relationship that appears in Table 2 is designated as one-to-many because there will be many records in the Groups table that have the same GroupID, but different ValueID, that serve to define the group. In Figure 1, the Group relationship is labeled 1..\* at the DataValues table and 0..\* at the Groups table. This indicates that a group may comprise one or more data values and that a data value may be included in 0 or more groups. Similarly, there will be many records in the DerivedFrom table that have the same DerivedFromID, but different ValueID that serve to define the group of data values from which a data value is derived. Logically a data value should not be in a DerivedFrom group upon which it is derived from. If this can be programmatically checked by the system, then this sort of circularity error could be prevented.

The method description in the Methods table associated with a data value that has a DerivedFromID should describe the method used for deriving the particular data value from other data values (e.g. calculating discharge from a number of velocity measurements across a stream). The relationship between the DataValues table DerivedFromID field and DerivedFrom table DerivedFromID field is many-to-many (\*<->\*) because it can occur that the same group of data values is used to derive more than one derived data value. In Figure 1, the AreDerivedFrom relationship between the data values and DerivedFrom table actually depicts both relationships between these tables listed in Table 2. The AreDerivedFrom relationship is labeled 1..\* at the DataValues table and 0..\* at the DerivedFrom table to indicate that a derived from group may comprise 1 or more data values and that a data value may be a member of 0 or more derived from groups.

### Dates and Times

Unambiguous interpretation of date and time information requires specification of the time zone or offset from universal time (UTC). A UTCOffset field is included in the DataValues table to ensure that local times recorded in the database can be referenced to standard time and to enable comparison of results across databases that may store data values collected in different time zones (e.g. compare data values from one hydrologic observatory to those collected at another hydrologic observatory located across the country). A design choice here was to have UTCOffset as a record level qualifier because even though the time zone, and hence offset, is likely the same for all measurements at a site, the offset may change due to daylight savings. Some investigators may run data loggers on UTC time, while others may use local time adjusting for daylight saving time. To avoid the necessity to keep track of the system used, or impose a system that might be cumbersome and lead to errors, we decided that if the offset was always recorded, the precise time would be unambiguous and would reduce the chance for interpretation errors. A field DateTimeUTC is also included as a record level attribute associated with each data value. This provides a consistent time for querying and sorting data values. There is a level of redundancy between LocalDateTime, UTCOffset and DateTimeUTC. Only two are required to calculate the third. For simplicity and clarity we retain all three. A specific database implementation may choose to retain only two and calculate the third on the fly. ODM data loaders should only require two of the quantities to be input and should then calculate the third.

The separation of the date and time specification into two variables, `LocalDateTime` and `UTCOffset`, in the generic conceptual model may be handled differently within specific implementations. In one specific implementation these may be grouped in one text field in standard (e.g. ISO 8601) format such as `YYYY-MM-DDhh:mm:ss.sss:UTCOffset` (e.g. `2006-03-2516:19:56.232:-7`), while in another format the date and time may be specified as the number of fractional days from an origin (e.g. Excel represents the above date as the following number `38801.6805` and allows the user to specify the format for display) with `UTCOffset` as a separate attribute. In general we expect specific implementations to take advantage of the representation of date time objects provided by the implementation software, but to expose the `LocalDateTime` and `UTCOffset` to users so that time may be unambiguously interpreted. In the `SeriesCatalog` table, begin and end times for each data series are represented by the attributes `BeginDateTime`, `EndDateTime`, `BeginDateTimeUTC`, and `EndDateTimeUTC`. The UTC offset may be derived from the difference between the UTC and local times. Because local time may change (e.g. with daylight savings) it is important during the derivation of the `SeriesCatalog` table that identification of the first and last records be based on UTC time and that local times be read from the corresponding records, rather than using a min or a max function on local times which can result in an error.

### Support Scale

In interpreting data values that comprise a time series it is important to know the scale information associated with the data values. Blöschl and Sivapalan (1995) review the important issues. Any set of data values is quantified by a scale triplet comprising support, spacing, and extent as illustrated in Figure 6.

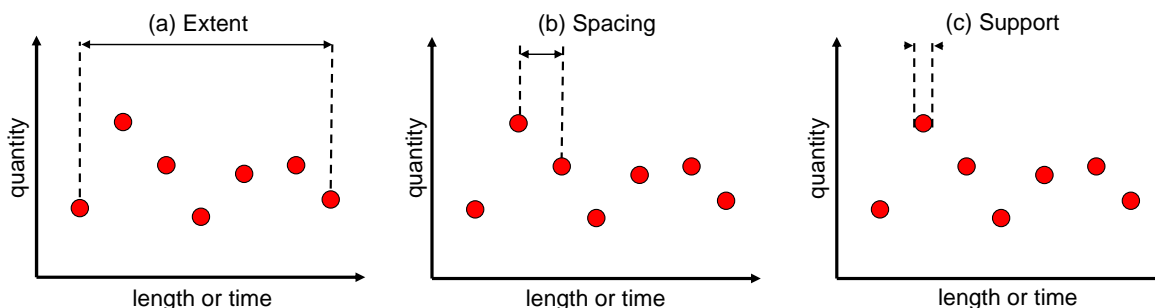


Figure 6. The scale triplet of measurements (a) extent, (b) spacing, (c) support (from Blöschl, 1996).

Extent is the full range over which the measurements occur, spacing is the spacing between measurements, and support is the averaging interval or footprint implicit in any measurement. In ODM, extent and spacing are properties of multiple measurements and are defined by the `LocalDateTime` or `DateTimeUTC` associated with data values. We have included a field called `TimeSupport` in the `Variables` table to explicitly quantify support. Figure 7 shows some of the implications associated with support, spacing, and extent in the interpretation of time series data values.

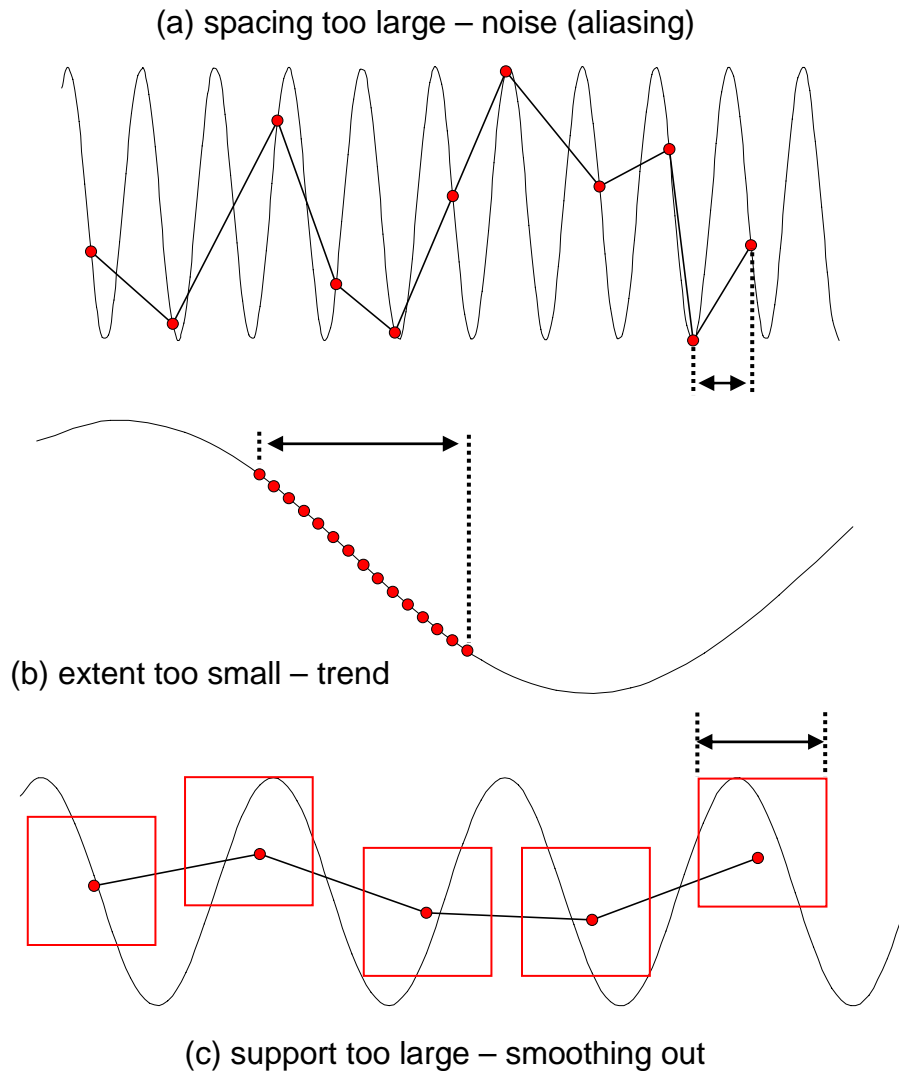


Figure 7. The effect of sampling for measurement scales not commensurate with the process scale: (a) spacing larger than the process scale causes aliasing in the data; (b) extent smaller than the process scale causes a trend in the data; (c) support larger than the process scale causes excessive smoothing in the data (adapted from Blöschl, 1996).

The concepts of scale described here apply in spatial as well as time dimensions. However, TimeSupport is only used to quantify support in the time dimension. The spatial support associated with a specific measurement method needs to be given or implied in the methods description in the Methods table. The next section indicates how time support should be specified for the different types of data.

### Data Types

In the ODM, the following data types are defined. These are specified by the DataType field in the Variables table.



1. *Continuous* data – the phenomenon, such as streamflow,  $Q(t)$  is specified at a particular instant in time and measured with sufficient frequency (small spacing) to be interpreted as a continuous record of the phenomenon. Time support may be specified as 0 if the measurements are instantaneous, or given a value that represents the time averaging inherent in the measurement method or device.
2. *Sporadic* data – the phenomenon is sampled at a particular instant in time but with a frequency that is too coarse for interpreting the record as continuous. This would be the case when the spacing is significantly larger than the support and the time scale of fluctuation of the phenomenon, such as for example infrequent water quality samples. As for continuous data, time support may be specified as 0 if the measurements are instantaneous, or given a value that represents the time averaging inherent in the measurement method or device.
3. *Cumulative* data – the data represents the cumulative value of a variable measured or calculated up to a given instant of time, such as cumulative volume of flow or cumulative

precipitation:  $V(t) = \int_0^t Q(\tau) d\tau$ , where  $\tau$  represents time in the integration over the

interval  $[0,t]$ . To unambiguously interpret cumulative data one needs to know the time origin. In the ODM we adopt the convention of using a cumulative record with a value of zero to initialize or reset cumulative data. With this convention, cumulative data should be interpreted as the accumulation over the time interval between the date and time of the zero record and the current record at the same site position. Site position is defined by a unique combination of SiteID, VariableID, OffsetValue and OffsetType. All four of these quantities comprise the unambiguous description of the position of an observation value and there may be multiple time series associated with multiple observation positions (e.g. redundant rain gauges with different offsets) at a location. The time support for a cumulative value should be specified as 0 if the measurement of the cumulative quantity is instantaneous, or given a value that represents the time averaging inherent in the measurement of the cumulative value at the end of the period of accumulation.

4. *Incremental* data – the data value represents the incremental value of a variable over a time interval  $\Delta t$  such as the incremental volume of flow, or incremental precipitation:

$\Delta V(t) = \int_t^{t+\Delta t} Q(\tau) d\tau$ . As for cumulative data, unambiguous interpretation requires

knowledge of the time increment. In the ODM we adopt the convention of using TimeSupport to specify the interval  $\Delta t$ , or the time interval to the next data value at the same position if TimeSupport is 0. This accommodates incremental type precipitation data that is only reported when the data value is non-zero, such as NCDC data. Such NCDC data is irregular, with the interpretation that precipitation is 0 if not reported unless qualifying comments designate otherwise. See example E.4 below for an illustration of how NCDC precipitation data is accommodated in the ODM.

5. *Average* data – the data value represents the average over a time interval, such as daily mean discharge or daily mean temperature:  $\bar{Q}(t) = \frac{\Delta V(t)}{\Delta t}$ . The averaging interval is quantified by TimeSupport in the case of regular data (as quantified by the IsRegular

field) and by the time interval from the previous data value at the same position for irregular data.

6. *Maximum* data – the data value is the maximum value occurring at some time during a time interval, such as annual maximum discharge or a daily maximum air temperature. Again unambiguous interpretation requires knowledge of the time interval. The ODM adopts the convention that the time interval is the TimeSupport for regular data and the time interval from the previous data value at the same position for irregular data.
7. *Minimum* data – the data value is the minimum value occurring at some time during a time interval, such as 7-day low flow for a year, or the daily minimum temperature. The time interval is defined similarly to Maximum data.
8. *Constant over interval* data – the data value is a quantity that can be interpreted as constant over the time interval to the next measurement.
9. *Categorical* data – the data value is a categorical rather than continuous valued quantity. Mapping from data values to categories is through the Categories table.

We anticipate that additional data types such as median, standard deviation, variance, and others may need to be added as users work with ODM.

#### Beginning of Interval Reporting Time for Interval Data Values

Data types 4 to 8 above apply to data values that occur over an interval of time. The date and time reported and entered in to the ODM database associated with each interval data value is the beginning time of the observation interval. This convention was adopted to be consistent with the way dates and times are represented in most common database management systems. It should be noted that using the beginning of the interval is not consistent with the time a data logger would log an observation value. Care should be exercised in adding data to the ODM to ensure that the beginning of interval convention is followed.

#### Time Series Data

A considerable portion of hydrologic observations data is in the form of time series. This was why the initial model was based on the Arc Hydro Time Series Data Model. The ODM design has not specifically highlighted time series capabilities; nevertheless, the data model has inherited the key components from the Arc Hydro Time Series Data Model to give it time series capability. In particular one variable DataType is “Continuous,” which is designed to indicate that the data values are collected with sufficient frequency as to be interpreted as a smooth time series. The IsRegular field also facilitates time series analysis because certain time series operations (e.g., Fourier Analysis) are predisposed to regularly sampled data. At first glance it may appear that there is redundancy between the IsRegular field and the DataType “Continuous,” but we chose to keep these separate because there are regularly sampled quantities for which it is not reasonable to interpret the data values as “Continuous.” For example, monthly grab samples of water quality are not continuous, but are better categorized as having DataType “Sporadic.” Note that ODM does not explicitly store the time interval between measurements, nor does it indicate where a continuous series has data gaps. Both of these are required for time series analysis, but are inherently not properties of single measurements. The time interval is the time difference between sequential regular measurements, something that can be easily computed from date and time values by analysis tools. The inference of measurement gaps (and

what to do about them) from date and time values we also regard as analysis functionality left for a Hydrologic Analysis System to handle.

### Categorical Variables

In ODM, categorical or ordinal variables are stored in the same table as continuous valued ‘real’ variables through a numerical encoding of the categorical data value as a ‘real’ data value. The Categories table then associates, for each variable, a data value with an associated category description. This is a somewhat cumbersome construct because real valued quantities are being used as database keys. We do not see this as a significant shortcoming though, because typically, in our judgment, only a small fraction of hydrologic observations will be categorical. The Categories table stores the categories associated with categorical data values. If a Variable has a DataType that is “Categorical” then the VariableID must match one or more VariableIDs in Categories that define the mapping between DataValues and Categories. The CategoryDescription field in the Categories table defines the category.

### Samples and Methods

At first glance there may appear to be redundancy between the information in the Samples table and Methods table. However, the samples table is intended to only be used where data values are derived from a physical sample that is later analyzed in a laboratory (e.g., a water chemistry sample or biological sample). The SampleID that links into the Samples table provides tracking of the specific physical sample used to derive each measurement and, by reference to information in the LabMethods table, the laboratory methods and protocols followed. The Methods table refers to the method of field data collection, which may specify “how” a physical observation was made or collected (e.g., from an automated sampler or collected manually), but is also used to specify the measurement method associated with an in-situ measurement instrument such as a weir, turbidity sensor, dissolved oxygen sensor, humidity sensor, or temperature sensor.

### Data Qualifiers

Each record in the DataValues table has an attribute called QualifierID that references the Qualifiers table. Each QualifierID in the Qualifiers table has attributes QualifierCode and QualifierDescription that provide qualifying information that can note anything unusual or problematic about individual observations such as, for example, "holding time for analysis exceeded" or "incomplete or inexact daily total." Specification of a QualifierID in the DataValues table is optional, with the inference that if a QualifierID is not specified then the corresponding data value is not qualified.

### Quality Control Level Encoding

Each data value in the DataValues table has an attribute called QualityControlLevelID that references the QualityControlLevels table and is designed to record the level of quality control processing that the data value has been subjected to at the level of data series. Quality control level is one of the attributes (together with site, variable, method, and source) used to uniquely identify data series. Each quality control level is uniquely identified by its QualityControlLevelID; however, each level also has a text QualityControlLevelCode that, along with a Definition and Explanation, provides a more descriptive encoding of the quality control level. The default quality control level system used by ODM applies integer values between 0

and 4 (converted to text strings) as the QualityControlLevelCodes. Other custom systems for QualityControlLevelCodes can be used (e.g., 0.1, 0.2 to represent raw data that is progressing through a quality control work sequence, or text strings such as “Raw” or “Processed”). The following 0 – 4 QualityControlLevelCode definitions are adapted from those used by other similar systems, such as NASA, Earthscope and Ameriflux (e.g. [http://ilrs.gsfc.nasa.gov/reports/ilrs\\_reports/9809\\_attach7a.html](http://ilrs.gsfc.nasa.gov/reports/ilrs_reports/9809_attach7a.html), <http://public.ornl.gov/ameriflux/available.shtml> accessed 3/6/2007) and are suggested so that CUAHSI ODM is consistent with the practice of other data systems:

- **QualityControlLevelCode = “0” - Raw Data**  
Raw data is defined as unprocessed data and data products that have not undergone quality control. Depending on the data type and data transmission system, raw data may be available within seconds or minutes after real-time. *Examples include real time precipitation, streamflow and water quality measurements.*
- **QualityControlLevelCode = “1” – Quality Controlled Data**  
Quality controlled data have passed quality assurance procedures such as routine estimation of timing and sensor calibration or visual inspection and removal of obvious errors. *An example is USGS published streamflow records following parsing through USGS quality control procedures.*
- **QualityControlLevelCode = “2” –Derived Products**  
Derived products require scientific and technical interpretation and include multiple-sensor data. *An example might be basin average precipitation derived from rain gages using an interpolation procedure.*
- **QualityControlLevelCode = “3” –Interpreted Products**  
These products require researcher (PI) driven analysis and interpretation, model-based interpretation using other data and/or strong prior assumptions. *An example is basin average precipitation derived from the combination of rain gages and radar return data.*
- **QualityControlLevelCode = “4” –Knowledge Products**  
These products require researcher (PI) driven scientific interpretation and multidisciplinary data integration and include model-based interpretation using other data and/or strong prior assumptions. *An example is percentages of old or new water in a hydrograph inferred from an isotope analysis.*

These definitions for quality control level are stored in the QualityControlLevels table. These definitions are recommended for use, but users can define their own quality control level system. The QualityControlLevels table is not a controlled vocabulary, but specification of a quality control level for each data value is required. Appendix B of this document provides a discussion of how to handle data versioning in terms of quality control levels (using the levels defined above), data series editing, and data series creation.

## Metadata

ODM has been designed to contain all the core elements of the CUAHHSI HIS metadata system (<http://www.cuahsi.org/his/metadata.html>) required for compliance with evolving standards such as the draft ISO 19115. In its design, the ODM embodies much record, variable, and site level metadata. Dataset and project level metadata required by these standards, such as TopicCategory, Title, and Abstract are included in a table called ISOMetaData linked to each data source.

## Reference Documents

The Methods, Sources, LabMethods and ISOMetaData tables contain fields that can be used to store links to source or reference information. At the general conceptual level of the ODM we do not specify how, or in what form these links to references or sources should be implemented. Options include using URLs or storing entire documents in the database. If external URLs are used it will be important as the database grows and is used over time to ensure that links or URLs included are stable. An alternative approach to external links is to exploit the capability of modern databases to store entire digital documents, such as an html or xml page, PDF document, or raw data file, within a field in the database. The capability therefore exists to instead have these links refer to a separate table that would actually contain this metadata information, instead of housing it in a separate digital library. There is some merit in this because then any data exported in ODM format could take with it the associated metadata required to completely define it as well as the raw data upon which it is derived. However, this has the disadvantage of increasing (perhaps substantially) the size of database file containing the data and being distributed to users.

## Controlled Vocabularies

The following tables in the ODM are tables where controlled vocabularies for the fields are required to maintain consistency and avoid the use of synonyms that can lead to ambiguity:

- CensorCodeCV
- DataTypeCV
- GeneralCategoryCV
- SampleMediumCV
- SampleTypeCV
- SpatialReferences
- SpeciationCV
- TopicCategoryCV
- Units
- ValueTypeCV
- VariableNameCV
- VerticalDatumCV

The initial contents of these controlled vocabularies are specified in the Microsoft SQL Server 2005 blank schema for the ODM. However, the ODM controlled vocabularies are dynamic. A central repository of current ODM controlled vocabulary terms is maintained on the ODM Website at <http://water.usu.edu/cuahsi/odm/>, together with the most recent version of the ODM

SQL Server 2005 blank schema, this design specifications document, and other tools for working with ODM. Users can submit new terms for the controlled vocabularies and can request changes to existing terms using functionality available on the ODM website (<http://water.usu.edu/cuahsi/odm/>). Functionality for updating local controlled vocabulary tables with new terms from the central ODM controlled vocabulary repository is provided in the ODM Tools software application, which is also available from the ODM website. The CUAHHSI HIS team welcomes input on the controlled vocabularies.

## **Examples**

The following examples show the capability of ODM to store different types of point observations. It is not possible in examples such as these to present all of the field values for all the tables. Because of this, the examples present selected fields and tables chosen to illustrate key capabilities of the data model. Refer to Appendix A for the complete definition of table and field contents.

### Streamflow - Gage Height and Discharge

Figure E.1 illustrates how both stream gage height measurements and the associated discharge estimates derived from the gage height measurements can be stored in the ODM. Note that gage height in feet and discharge in cubic feet per second are both in the same data table but with different VariableIDs that reference the Variables table, which specifies the VariableName, Units, and other quantities associated with these data values. The link between VariableID in the DataValues table and Variables table is shown. In this example, discharge measurements are derived from gage height (stage) measurements through a rating curve. The MethodID associated with each discharge record references into the Methods table that describes this and provides a URL that should contain metadata details for this method. The DerivedFromID in the DataValues table references into the DerivedFrom table that references back to the corresponding gage height in the DataValues table from which the discharge was derived.

**DataValues : Table**

ValueID	DataValue	ValueAccuracy	LocalDateTime	UTCOffset	SiteID	VariableID	MethodID	DerivedFromID
1	4.18		05/01/2006 0:00:00.000	-7	1	1	1	
97	748		05/01/2006 0:00:00.000	-7	1	2	1	
193	722	22.89831642	05/01/2006 0:00:00.000	-7	1	3	100	
2	4.18		05/01/2006 0:15:00.000	-7	1	1	1	
98	748		05/01/2006 0:15:00.000	-7	1	2	2	
3	4.17		05/01/2006 0:30:00.000	-7	1	1	1	
99	742		05/01/2006 0:30:00.000	-7	1	2	3	
4	4.17		05/01/2006 0:45:00.000	-7	1	1	1	
100	742		05/01/2006 0:45:00.000	-7	1	2	4	
5	4.17		05/01/2006 1:00:00.000	-7	1	1	1	
101	742		05/01/2006 1:00:00.000	-7	1	2	5	
6	4.17		05/01/2006 1:15:00.000	-7	1	1	1	
102	742		05/01/2006 1:15:00.000	-7	1	2	6	

**DerivedFrom : Table**

DerivedFromID	ValueID
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17

**Variables : Table**

VariableID	VariableCode	VariableName	VariableUnitsID	SampleMedium	ValueType	IsRegular	TimeSupport	TimeUnitsID	DataType	GeneralCategory	NoDataValue
1	00065	Gage height	1	Water	Field Observation	<input checked="" type="checkbox"/>	15	5	Continuous	Hydrologic	-9999
2	00060	Discharge	2	Water	Derived Value	<input checked="" type="checkbox"/>	15	5	Continuous	Hydrologic	-9999
3	00060	Discharge, daily average	2	Water	Derived Value	<input checked="" type="checkbox"/>	24	6	Average	Hydrologic	-9999
4	00300	Dissolved oxygen concentration	3	Water	Field Observation	<input type="checkbox"/>	0		Instantaneous	Water Quality	-9999

**Units : Table**

UnitsID	UnitsName	UnitsType	UnitsAbbreviation
1	Feet	Length	ft
2	Cubic feet per second	Flow	ft <sup>3</sup> /s
3	Milligrams per liter	Concentration	mg/L
4	Meters	Length	m
5	Minutes	Time	min
6	Hours	Time	hr

**Methods : Table**

MethodID	MethodDescription
1	Gage height measured with continuous data logger
2	Discharge derived from water stage using site specific rating curve
3	Daily average discharge derived from 15 minute continuous discharge values
4	Dissolved oxygen measured with a Hydrolab multiprobe field instrument

Figure E.1. Excerpts from tables illustrating the population of ODM with streamflow gage height (stage) and discharge data.

### Streamflow - Daily Average Discharge

Daily average streamflow is reported as an average of continuous 15 minute interval data values. Figure E.2 shows excerpts from tables illustrating the population of ODM with both the continuous discharge values and derived daily averages. The record giving the single daily average discharge with a value of 722 ft<sup>3</sup>/s in the DataValues table has a DerivedFromID of 100. This refers to multiple records in the DerivedFrom table, with associated ValueIDs 97, 98, 99, ... 113 shown. These refer to the specific 15 minute discharge values in the DataValues table used to derive the average daily discharge. VariableID in the DataValues table identifies the appropriate record in the Variables table specifying that this is a daily average discharge with units of ft<sup>3</sup>/s from UnitsID referencing in to the Units table. MethodID in the DataValues table identifies the appropriate record in the Methods table specifying that the method used to obtain this data value was daily averaging.

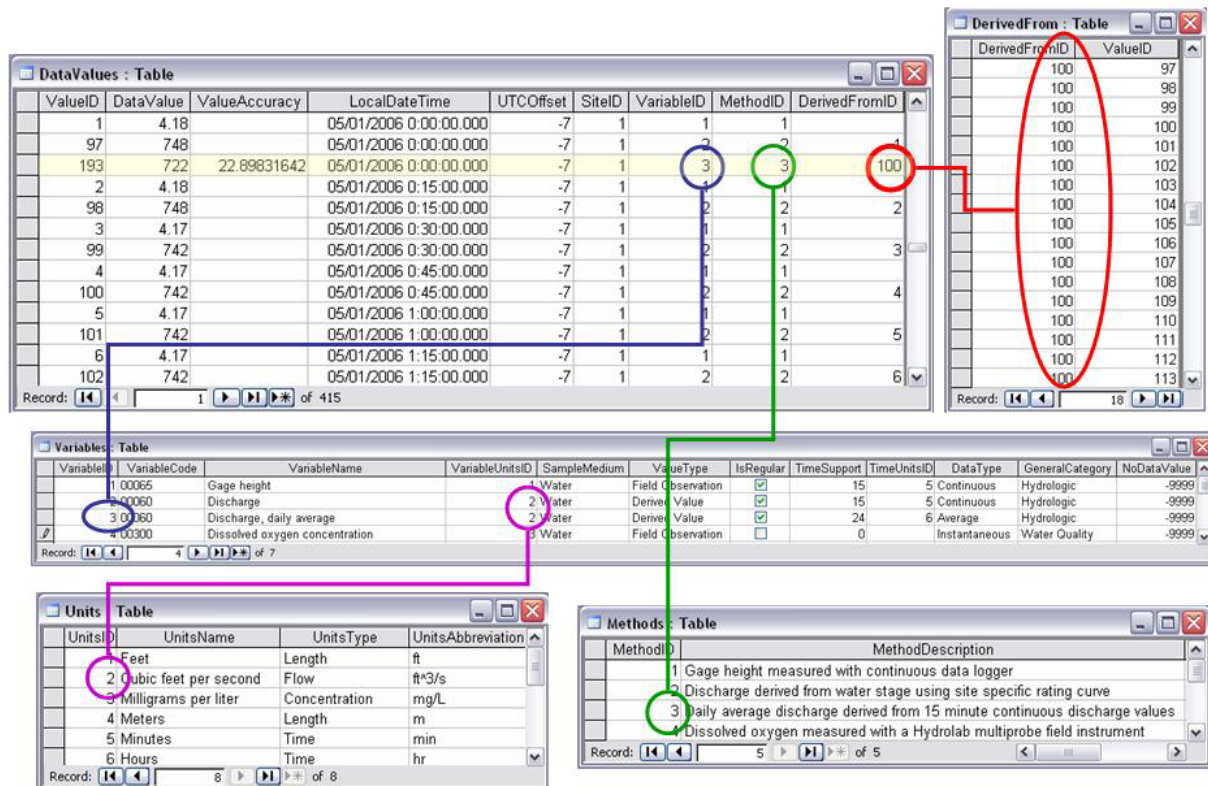


Figure E.2. Excerpts from tables illustrating the population of ODM with daily average discharge derived from 15 minute discharge values.

### Water Chemistry from a Profile in a Lake

Reservoir profile measurements provide an example of the logical grouping of data values and data values that have an offset in relationship to the location of the monitoring site. These measurements may be made simultaneously (by multiple instruments in the water column) or over a short time period (one instrument that is lowered from top to bottom). Figure E.3 shows an example of how these data would be stored in ODM. The OffsetTypes table and OffsetValue attribute is used to quantify the depth offset associated with each measurement. Each of the data values shown has an OffsetTypeID that references into the OffsetTypes table. The OffsetTypes table indicates that for this OffsetType the offset is “Depth below water surface.” The OffsetTypes table references into the Units table indicating that the OffsetUnits are meters, so OffsetValue in the DataValues table is in units of meters depth below the water surface. Each of the data values shown also has a VariableID that in the Variables table indicates that the variable measured was dissolved oxygen concentration in units of mg/L. Each of the data values shown also has a MethodID that in the Methods table indicates that dissolved oxygen was measured with a Hydrolab multiprobe. The data values shown are part of a logical group of data values representing the water chemistry profile in a lake. This is represented using the Groups table and GroupDescriptions table. The Groups table associates GroupID 1 with each of the ValueIDs of the data values belonging to the group. A description of this group is given in the GroupDescriptions table.



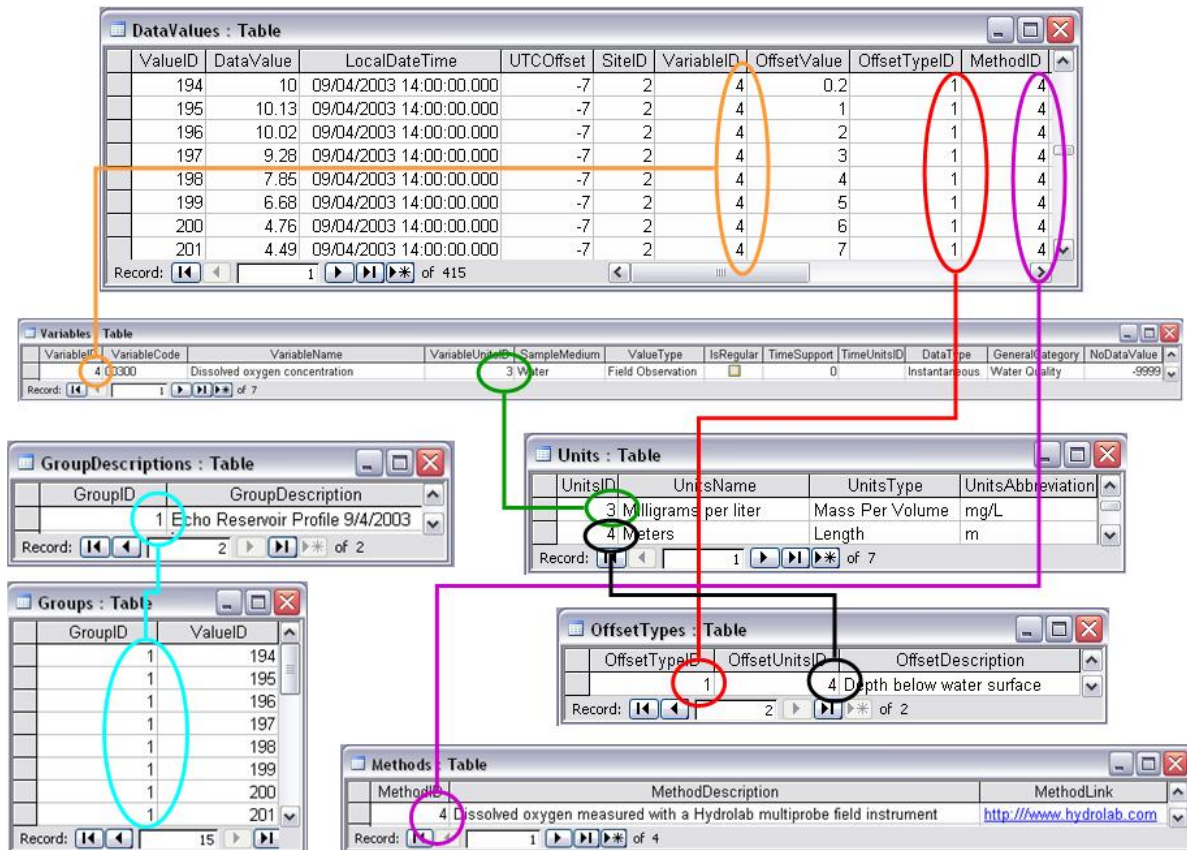


Figure E.3. Excerpts from tables illustrating the population of ODM with water chemistry data.

### NCDC Precipitation Data

Figure E.4 illustrates the representation of NCDC 15 minute precipitation data by ODM. The data includes 15 minute incremental data values as well as daily totals. Separate records in the Variables table are used for the 15 minute or daily total values. These data are reported at irregular intervals and only logged for time periods for which precipitation is non zero. This is accommodated by setting the IsRegular attribute associated with the variable to “False” and specifying the TimeSupport value as 15 or 24 and the TimeUnits as “Minutes” or “Hours”. The DataType of “Incremental” is used to indicate that these are incremental data values defined over the TimeSupport interval. The convention for incremental data (see above) is that when the time support is specified, it specifies the increment for irregular incremental data. When time support is specified as 0 it means the increment is from the previous data value at the same site position. Data qualifiers indicate periods where the data is missing. The method associated with each precipitation variable documents the convention that zero precipitation periods are not logged in this data acquired from NCDC. A data qualifier is also used to flag days where the precipitation total is incomplete due to the record being missing during part of the day.

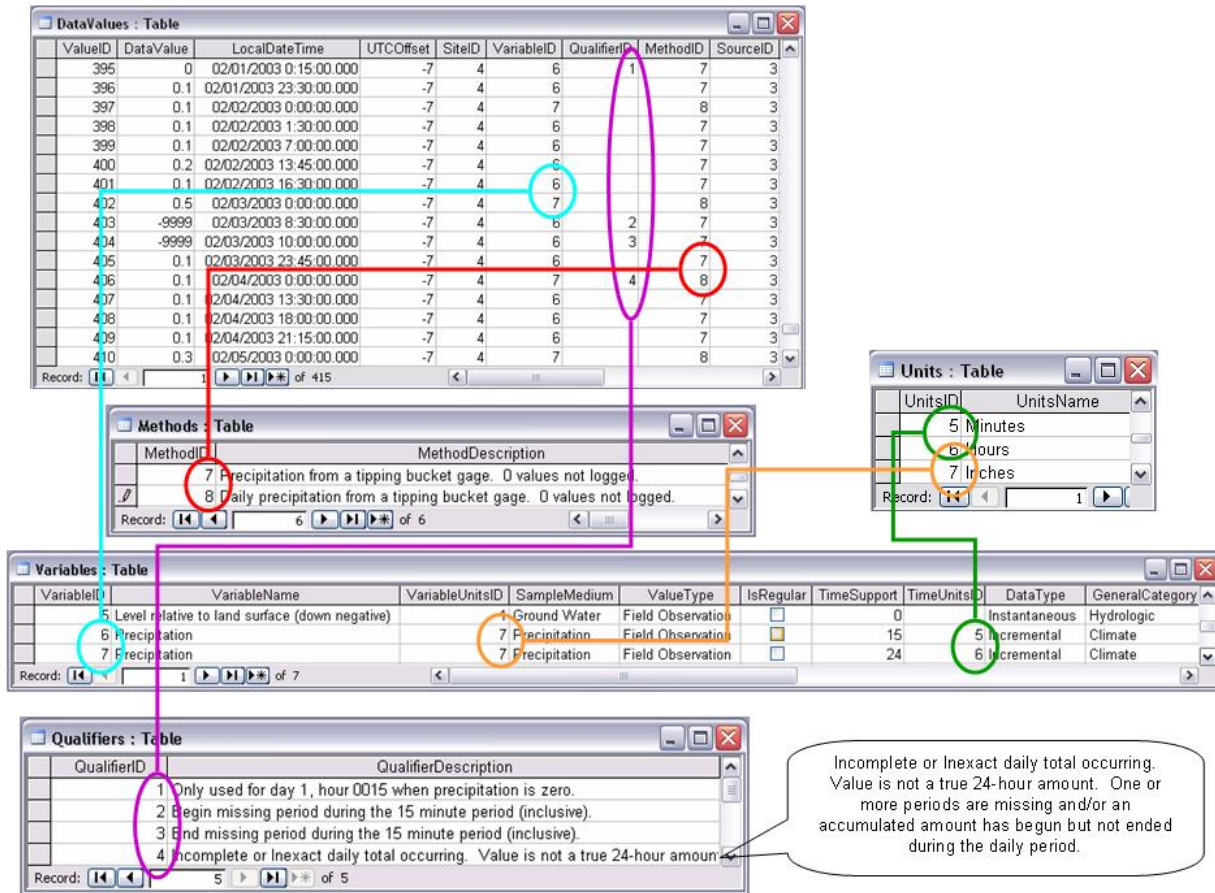


Figure E.4. Excerpts from tables illustrating the population of the ODM with NCDC Precipitation Data.

### Groundwater Level Data

The following is an example of how groundwater level data can be stored in ODM. In this example, the data values are the water table level relative to the ground surface reported as negative values. This example shows multiple data values of a single variable at a single site made by a single source that have been quality controlled as indicated by the QualityControlLevelID field in the QualityControlLevels table. The SiteID field in the DataValues table indicates the site in the Sites table that gives the location information about the monitoring site. In this case, the elevation is with respect to the NGVD29 datum as indicated in the VerticalDatum field, and latitude and longitude are with respect to the NAD27 datum as indicated in the SpatialReferences table. The VariableID field in the DataValues table references the appropriate record in the Variables table indicating information about the variable. The SourceID field in the DataValues table references the appropriate record in the Sources table giving information about the source of the data.

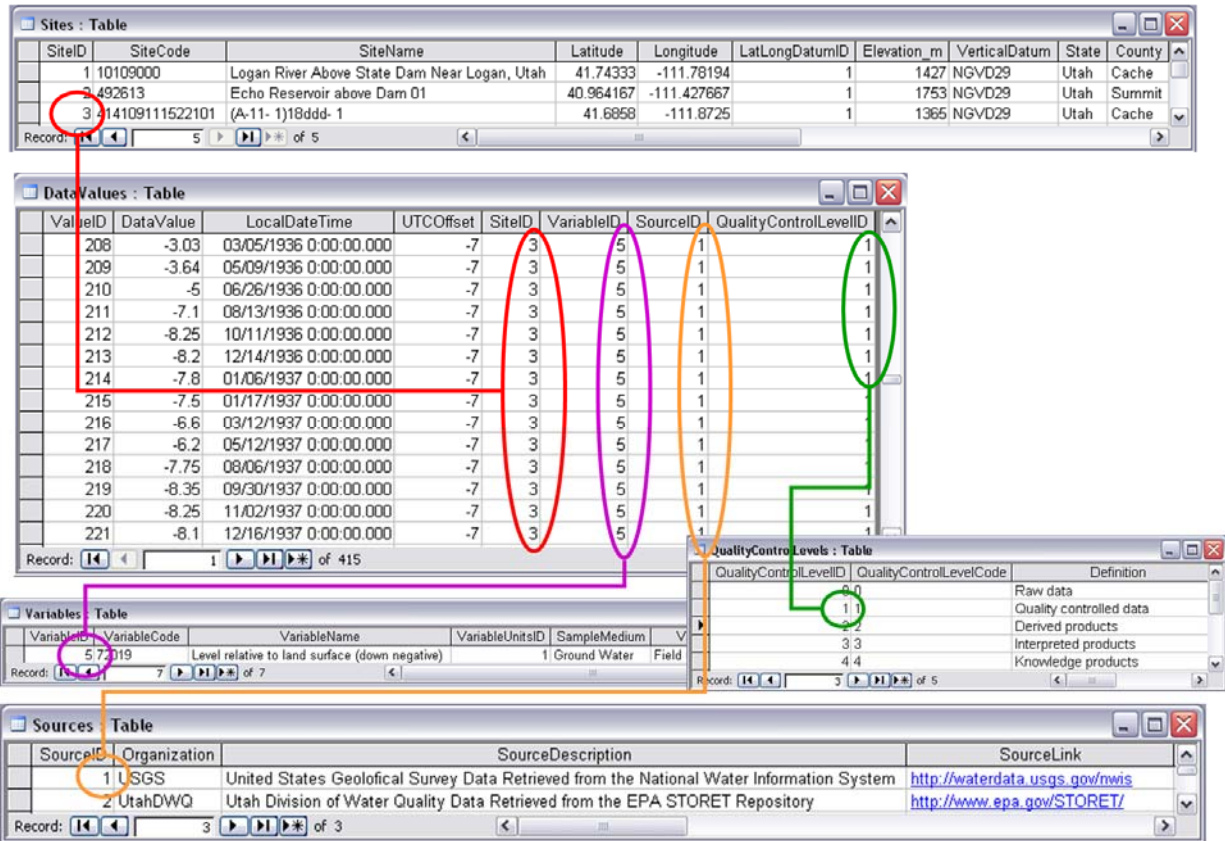


Figure E.5. Excerpts from tables illustrating the population of the ODM with irregularly sampled groundwater level data.

### Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant Nos. EAR 0412975 and 0413265. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation (NSF).

## References

- AIAA, (1995), Assessment of Wind Tunnel Data Uncertainty, American Institute of Aeronautics and Astronautics: AIAA S-071-1995.
- Blöschl, G., (1996), Scale and Scaling in Hydrology, Habilitationsschrift, Weiner Mitteilungen Wasser Abwasser Gewässer, Wien, 346 p.
- Blöschl, G. and M. Sivapalan, (1995), "Scale Issues in Hydrological Modelling: A Review," Hydrological Processes, 9(1995): 251-290.
- Horsburgh, J. S., D. G. Tarboton and D. R. Maidment, (2005), "A Community Data Model for Hydrologic Observations, Chapter 6," in Hydrologic Information System Status Report, Version 1, Edited by D. R. Maidment, p.102-135, <http://www.cuahsi.org/his/docs/HISStatusSept15.pdf>.
- Horsburgh, J. S., D. G. Tarboton, D. R. Maidment, and I. Zaslavsky, (2008), A relational model for environmental and water resources data, Water Resources Research, Vol. 44, W05406, doi:10.1029/2007WR006392.
- Maidment, D. R., ed. (2002), Arc Hydro GIS for Water Resources, ESRI Press, Redlands, CA, 203 p.
- Maidment, D. R., (2005), "A Data Model for Hydrologic Observations." Paper prepared for presentation at the CUAHSI Hydrologic Information Systems Symposium, University of Texas at Austin. March 7, 2005.
- Tarboton, D. G., (2005), "Review of Proposed CUAHSI Hydrologic Information System Hydrologic Observations Data Model." Utah State University. May 5, 2005.

## Appendix A. Observations Data Model Table and Field Structure

The following is a description of the tables in the observations data model, a listing of the fields contained in each table, a description of the data contained in each field and its data type, examples of the information to be stored in each field where appropriate, specific constraints imposed on each field, and discussion on how each field should be populated. Values in the example column should not be considered to be inclusive of all potential values, especially in the case of fields that require a controlled vocabulary. We anticipate that these controlled vocabularies will need to be extended and adjusted. Tables appear in alphabetical order.

Each table below includes a “Constraint” column. The value in this column designates each field in the table as one of the following:

Mandatory (M) – A value in this field is mandatory and cannot be NULL.

Optional (O) – A value in this field is optional and can be NULL.

Programmatically derived (P) – Inherits from the source field. The value in this field should be automatically populated as the result of a query and is not required to be input by the user.

Additional constraints are documented where appropriate in the Constraint column. In addition, where appropriate, each table contains a “Default Value” column. The value in this column is the default value for the associated field. The default value specifies the convention that should be followed when a value for the field is not specified. Below each table is a discussion of the rules and best practices that should be used in populating each table within ODM.

### Table: Categories

The Categories table defines the categories for categorical variables. Records are required for variables where DataType is specified as "Categorical." Multiple entries for each VariableID, with different DataValues provide the mapping from DataValue to category description.

Field Name	DataType	Description	Examples	Constraint
VariableID	Integer	Integer identifier that references the Variables record of a categorical variable.	45	M Foreign key
DataValue	Real	Numeric value that defines the category	1.0	M
CategoryDescription	Text (Unlimited)	Definition of categorical variable value	“Cloudy”	M

The following rules and best practices should be used in populating this table:

1. Although all of the fields in this table are mandatory, they need only be populated if categorical data are entered into the database. If there are no categorical data in the DataValues table, this table will be empty.
2. This table should be populated before categorical data values are added to the DataValues table.

**Table: CensorCodeCV**

The CensorCodeCV table contains the controlled vocabulary for censor codes. Only values from the Term field in this table can be used to populate the CensorCode field of the DataValues table.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for CensorCode.	“lt”, “gt”, “nc”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of CensorCode controlled vocabulary term. The definition is optional if the term is self explanatory.	“less than”, “greater than”, “not censored”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: DataTypeCV**

The DataTypeCV table contains the controlled vocabulary for data types. Only values from the Term field in this table can be used to populate the DataType field in the Variables table.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for DataType.	“Continuous”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of DataType controlled vocabulary term. The definition is optional if the term is self explanatory.	“A quantity specified at a particular instant in time measured with sufficient frequency (small spacing) to be interpreted as a continuous record of the phenomenon.”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: DataValues**

The DataValues table contains the actual data values.

Field Name	Data Type	Description	Example	Constraint	Default Value
ValueID	Integer Identity	Unique integer identifier for each data value.	43	M Unique Primary key	
DataValue	Real	The numeric value of the observation. For Categorical variables, a number is stored here. The Variables table has DataType as Categorical and the Categories table maps from the DataValue onto Category Description.	34.5	M	
ValueAccuracy	Real	Numeric value that describes the measurement accuracy of the data value. If not given, it is interpreted as unknown.	4	O	NULL
LocalDateTime	Date/Time	Local date and time at which the data value was observed. Represented in an implementation specific format.	9/4/2003 7:00:00 AM	M	
UTCOffset	Real	Offset in hours from UTC time of the corresponding LocalDateTime value.	-7	M	
DateTimeUTC	Date/Time	Universal UTC date and time at which the data value was observed. Represented in an implementation specific format.	9/4/2003 2:00:00 PM	M	
SiteID	Integer	Integer identifier that references the site at which the observation was measured. This links data values to their locations in the Sites table.	3	M Foreign key	
VariableID	Integer	Integer identifier that references the variable that was measured. This links data values to their variable in the Variables table.	5	M Foreign key	
OffsetValue	Real	Distance from a datum or control point to the point at which a data value was observed. If not given the OffsetValue is inferred to be 0, or not relevant/necessary.	2.1	O	NULL = No Offset
OffsetTypeID	Integer	Integer identifier that references the measurement offset type in the OffsetTypes table.	3	O Foreign key	NULL = No Offset
CensorCode	Text (50)	Text indication of whether the data value is censored from the CensorCodeCV controlled vocabulary.	“nc”	M Foreign key	“nc” = Not Censored

Field Name	Data Type	Description	Example	Constraint	Default Value
QualifierID	Integer	Integer identifier that references the Qualifiers table. If Null, the data value is inferred to not be qualified.	4	O Foreign key	NULL
MethodID	Integer	Integer identifier that references method used to generate the data value in the Methods table.	3	M Foreign key	0 = No method specified
SourceID	Integer	Integer identifier that references the record in the Sources table giving the source of the data value.	5	M Foreign key	
SampleID	Integer	Integer identifier that references into the Samples table. This is required only if the data value resulted from a physical sample processed in a lab.	7	O Foreign key	NULL
DerivedFromID	Integer	Integer identifier for the derived from group of data values that the current data value is derived from. This refers to a group of derived from records in the DerivedFrom table. If NULL, the data value is inferred to not be derived from another data value.	5	O	NULL
QualityControlLevelID	Integer	Integer identifier giving the level of quality control that the value has been subjected to. This references the QualityControlLevels table.	1	M Foreign key	-9999 = Unknown

The following rules and best practices should be used in populating this table:

1. ValueID is the primary key, is mandatory, and cannot be NULL. This field should be implemented as an autonumber/identity field. When data values are added to this table, a unique integer ValueID should be assigned to each data value by the database software such that the primary key constraint is not violated.
2. Each record in this table must be unique. This is enforced by a unique constraint across all of the fields in this table (excluding ValueID) so that duplicate records are avoided.
3. The LocalDateTime, UTCOffset, and DateTimeUTC must all be populated. Care must be taken to ensure that the correct UTCOffset is used, especially in areas that observe daylight saving time. If LocalDateTime and DateTimeUTC are given, the UTCOffset can be calculated as the difference between the two dates. If LocalDateTime and UTCOffset are given, DateTimeUTC can be calculated.
4. SiteID must correspond to a valid SiteID from the Sites table. When adding data for a new site to the ODM, the Sites table should be populated prior to adding data values to the DataValues table.
5. VariableID must correspond to a valid VariableID from the Variables table. When adding data for a new variable to the ODM, the Variables table should be populated prior to adding data values for the new variable to the DataValues table.
6. OffsetValue and OffsetTypeID are optional because not all data values have an offset. Where no offset is used, both of these fields should be set to NULL indicating that the data values do not have an offset. Where an OffsetValue is specified, an OffsetTypeID



must also be specified and it must refer to a valid OffsetTypeID in the OffsetTypes table. The OffsetTypes table should be populated prior to adding data values with a particular OffsetTypeID to the DataValues table.

7. CensorCode is mandatory and cannot be NULL. A default value of “nc” is used for this field. Only Terms from the CensorCodeCV table should be used to populate this field.
8. The QualifierID field is optional because not all data values have qualifiers. Where no qualifier applies, this field should be set to NULL. When a QualifierID is specified in this field it must refer to a valid QualifierID in the Qualifiers table. The Qualifiers table should be populated prior to adding data values with a particular QualifierID to the DataValues Table.
9. MethodID must correspond to a valid MethodID from the Methods table and cannot be NULL. A default value of 0 is used in the case where no method is specified or the method used to create the observation is unknown. The Methods table should be populated prior to adding data values with a particular MethodID to the DataValues table.
10. SourceID must correspond to a valid SourceID from the Sources table and cannot be NULL. The Sources table should be populated prior to adding data values with a particular SourceID to the DataValues table.
11. SampleID is optional and should only be populated if the data value was generated from a physical sample that was sent to a laboratory for analysis. The SampleID must correspond to a valid SampleID in the Samples table, and the Samples table should be populated prior to adding data values with a particular SampleID to the DataValues table.
12. DerivedFromID is optional and should only be populated if the data value was derived from other data values that are also stored in the ODM database.
13. QualityControlLevelID is mandatory, cannot be NULL, and must correspond to a valid QualityControlLevelID in the QualityControlLevels table. A default value of -9999 is used for this field in the event that the QualityControlLevelID is unknown. The QualityControlLevels table should be populated prior to adding data values with a particular QualityControlLevelID to the DataValues table.

**Table: DerivedFrom**

The DerivedFrom table contains the linkage between derived data values and the data values that they were derived from.

Field Name	Data Type	Description	Examples	Constraint
DerivedFromID	Integer	Integer identifying the group of data values from which a quantity is derived.	3	M
ValueID	Integer	Integer identifier referencing data values that comprise a group from which a quantity is derived. This corresponds to ValueID in the DataValues table.	1,2,3,4,5	M

The following rules and best practices should be used in populating this table:

1. Although all of the fields in this table are mandatory, they need only be populated if derived data values and the data values that they were derived from are entered into the database. If there are no derived data in the DataValues table, this table will be empty.

**Table: GeneralCategoryCV**

The GeneralCategoryCV table contains the controlled vocabulary for the general categories associated with Variables. The GeneralCategory field in the Variables table can only be populated with values from the Term field of this controlled vocabulary table.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for GeneralCategory.	“Hydrology”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of GeneralCategory controlled vocabulary term. The definition is optional if the term is self explanatory.	“Data associated with hydrologic variables or processes.”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: GroupDescriptions**

The GroupDescriptions table lists the descriptions for each of the groups of data values that have been formed.

Field Name	Data Type	Description	Example	Constraint
GroupID	Integer Identity	Unique integer identifier for each group of data values that has been formed. This also references to GroupID in the Groups table.	4	M Unique Primary key
GroupDescription	Text (unlimited)	Text description of the group.	“Echo Reservoir Profile 7/7/2005”	O

The following rules and best practices should be used in populating this table:

1. This table will only be populated if groups of data values have been created in the ODM database.
2. The GroupID field is the primary key, must be a unique integer, and cannot be NULL. It should be implemented as an auto number/identity field.
3. The GroupDescription can be any text string that describes the group of observations.

**Table: Groups**

The Groups table lists the groups of data values that have been created and the data values that are within each group.

Field Name	Data Type	Description	Example	Constraint
GroupID	Integer	Integer ID for each group of data values that has been formed.	4	M Foreign key
ValueID	Integer	Integer identifier for each data value that belongs to a group. This corresponds to ValueID in the DataValues table	2,3,4	M Foreign key

The following rules and best practices should be used in populating this table:

1. This table will only be populated if groups of data values have been created in the ODM database.
2. The GroupID field must reference a valid GroupID from the GroupDescriptions table, and the GroupDescriptions table should be populated for a group prior to populating the Groups table.

**Table: ISOMetadata**

The ISOMetadata table contains dataset and project level metadata required by the CUAHSI HIS metadata system (<http://www.cuahsi.org/his/documentation.html>) for compliance with standards such as the draft ISO 19115 or ISO 8601. The mandatory fields in this table must be populated to provide a complete set of ISO compliant metadata in the database.

Field Name	Data Type	Description	Example	Constraint	Default Value
MetadataID	Integer Identity	Unique integer ID for each metadata record.	4	M Unique Primary key	
TopicCategory	Text (255)	Topic category keyword that gives the broad ISO19115 metadata topic category for data from this source. The controlled vocabulary of topic category keywords is given in the TopicCategoryCV table.	“inlandWaters”	M Foreign key	“Unknown”
Title	Text (255)	Title of data from a specific data source.		M Cannot contain tab, line feed, or carriage return characters	“Unknown”
Abstract	Text (unlimited)	Abstract of data from a specific data source.		M	“Unknown”

Field Name	Data Type	Description	Example	Constraint	Default Value
ProfileVersion	Text (255)	Name of metadata profile used by the data source	“ISO8601”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
MetadataLink	Text (500)	Link to additional metadata reference material.		O	NULL

The following rules and best practices should be used in populating this table:

1. The MetadataID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
2. All of the fields in this table are mandatory and cannot be NULL except for the MetadataLink field.
3. The TopicCategory field should only be populated with terms from the TopicCategoryCV table. The default controlled vocabulary term is “Unknown.”
4. The Title field should be populated with a brief text description of what the referenced data represent. This field can be populated with “Unknown” if there is no title for the data.
5. The Abstract field should be populated with a more complete text description of the data that the metadata record references. This field can be populated with “Unknown” if there is no abstract for the data.
6. The ProfileVersion field should be populated with the version of the ISO metadata profile that is being used. This field can be populated with “Unknown” if there is no profile version for the data.
7. One record with a MetadataID = 0 should exist in this table with TopicCategory, Title, Abstract, and ProfileVersion = “Unknown” and MetadataLink = NULL. This record should be the default value for sources with unknown/unspecified metadata.

### Table: LabMethods

The LabMethods table contains descriptions of the laboratory methods used to analyze physical samples for specific constituents.

Field Name	Data Type	Description	Example	Constraint	Default Value
LabMethodID	Integer Identity	Unique integer identifier for each laboratory method. This is the key used by the Samples table to reference a laboratory method.	6	M Unique Primary key	

LabName	Text (255)	Name of the laboratory responsible for processing the sample.	“USGS Atlanta Field Office”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
LabOrganization	Text (255)	Organization responsible for sample analysis.	“USGS”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
LabMethodName	Text (255)	Name of the method and protocols used for sample analysis.	“USEPA-365.1”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
LabMethodDescription	Text (unlimited)	Description of the method and protocols used for sample analysis.	“Processed through Model *** Mass Spectrometer”	M	“Unknown”
LabMethodLink	Text (500)	Link to additional reference material on the analysis method.		O	NULL

The following rules and best practices should be used when populating this table:

1. The LabMethodID field is the primary key, must be a unique integer, and cannot be NULL. It should be implemented as an auto number/identity field.
2. All of the fields in this table are required and cannot be null except for the LabMethodLink.
3. The default value for all of the required fields except for the LabMethodID is “Unknown.”
4. A single record should exist in this table where the LabMethodID = 0 and the LabName, LabOrganization, LabMethodName, and LabMethodDescription fields are equal to “Unknown” and the LabMethodLink = NULL. This record should be used to identify samples in the Samples table for which nothing is known about the laboratory method used to analyze the sample.

**Table: Methods**

The Methods table lists the methods used to collect the data and any additional information about the method.

Field Name	Data Type	Description	Example	Constraint	Default Value
MethodID	Integer Identity	Unique integer ID for each method.	5	M Unique Primary key	
MethodDescription	Text (unlimited)	Text description of each method.	“Specific conductance measured using a Hydrolab” or “Streamflow measured using a V notch weir with dimensions xxx”	M	
MethodLink	Text (500)	Link to additional reference material on the method.		O	NULL

The following rules and best practices should be used when populating this table:

1. The MethodID field is the primary key, must be a unique integer, and cannot be NULL.
2. There is no default value for the MethodDescription field in this table. Rather, this table should contain a record with MethodID = 0, MethodDescription = “Unknown”, and MethodLink = NULL. A MethodID of 0 should be used as the MethodID for any data values for which the method used to create the value is unknown (i.e., the default value for the MethodID field in the DataValues table is 0).
3. Methods should describe the manner in which the observation was collected (i.e., collected manually, or collected using an automated sampler) or measured (i.e., measured using a temperature sensor or measured using a turbidity sensor). Details about the specific sensor models and manufacturers can be included in the MethodDescription.

**Table: ODM Version**

The ODM Version table has a single record that records the version of the ODM database. This table must contain a valid ODM version number. This table will be pre-populated and should not be edited.

Field Name	Data Type	Description	Example	Constraint
VersionNumber	Text (50)	String that lists the version of the ODM database.	“1.1”	M Cannot contain tab, line feed, or carriage return characters

**Table: OffsetTypes**

The OffsetTypes table lists full descriptive information for each of the measurement offsets.

Field Name	Data Type	Description	Example	Constraint
OffsetTypeID	Integer Identity	Unique integer identifier that identifies the type of measurement offset.	2	M Unique Primary key
OffsetUnitsID	Integer	Integer identifier that references the record in the Units table giving the Units of the OffsetValue.	1	M Foreign key
OffsetDescription	Text (unlimited)	Full text description of the offset type.	“Below water surface” “Above Ground Level”	M

The following rules and best practices should be followed when populating this table:

1. Although all three fields in this table are mandatory, this table will only be populated if data values measured at an offset have been entered into the ODM database.
2. The OffsetTypeID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
3. The OffsetUnitsID field should reference a valid ID from the UnitsID field in the Units table. Because the Units table is a controlled vocabulary, only units that already exist in the Units table can be used as the units of the offset.
4. The OffsetDescription field should be filled in with a complete text description of the offset that provides enough information to interpret the type of offset being used. For example, “Distance from stream bank” is ambiguous because it is not known which bank is being referred to.

**Table: Qualifiers**

The Qualifiers table contains data qualifying comments that accompany the data.

Field Name	Data Type	Description	Example	Constraint	Default Value
QualifierID	Integer Identity	Unique integer identifying the data qualifier.	3	M Unique Primary key	
QualifierCode	Text (50)	Text code used by organization that collects the data.	“e” (for estimated) or “a” (for approved) or “p” (for provisional)	O Cannot contain space, tab, line feed, or carriage return characters	NULL
QualifierDescription	Text (unlimited)	Text of the data qualifying comment.	“Holding time for sample analysis exceeded”	M	

This table will only be populated if data values that have data qualifying comments have been added to the ODM database. The following rules and best practices should be used when populating this table:

1. The QualifierID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.

**Table: QualityControlLevels**

The QualityControlLevels table contains the quality control levels that are used for versioning data within the database.

Field Name	Data Type	Description	Example	Constraint
QualityControlLevelID	Integer Identity	Unique integer identifying the quality control level.	0, 1, 2, 3, 4, 5	M Unique Primary key
QualityControlLevelCode	Text (50)	Code used to identify the level of quality control to which data values have been subjected.	“1”, “1.1”, “Raw”, “QC Checked”	M Cannot contain tab, line feed, or carriage return characters
Definition	Text (255)	Definition of Quality Control Level.	“Raw Data”, “Quality Controlled Data”	M Cannot contain tab, line feed, or carriage return characters
Explanation	Text (unlimited)	Explanation of Quality Control Level	“Raw data is defined as unprocessed data and data products that have not undergone quality control.”	M

This table is pre-populated with quality control levels 0 through 4 within the ODM. The following rules and best practices should be used when populating this table:

1. The QualityControlLevelID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
2. It is suggested that the pre-populated system of quality control level codes (i.e., QualityControlLevelCodes 0 – 4) be used. If the pre-populated list is not sufficient, new quality control levels can be defined. A quality control level code of -9999 is suggested for data whose quality control level is unknown.



**Table: SampleMediumCV**

The SampleMediumCV table contains the controlled vocabulary for sample media.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for sample media.	“Surface Water”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of sample media controlled vocabulary term. The definition is optional if the term is self explanatory.	“Sample taken from surface water such as a stream, river, lake, pond, reservoir, ocean, etc.”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: Samples**

The Samples table gives information about physical samples analyzed in a laboratory.

Field Name	Data Type	Description	Example	Constraint	Default Value
SampleID	Integer Identity	Unique integer identifier that identifies each physical sample.	3	M Unique Primary key	
SampleType	Text (255)	Controlled vocabulary specifying the sample type from the SampleTypeCV table.	“FD”, “PB”, “SW”, “Grab Sample”	M Foreign key	“Unknown”
LabSampleCode	Text (50)	Code or label used to identify and track lab sample or sample container (e.g. bottle) during lab analysis.	“AB-123”	M Unique Cannot contain tab, line feed, or carriage return characters	
LabMethodID	Integer	Unique identifier for the laboratory method used to process the sample. This references the LabMethods table.	4	M Foreign key	0 = Nothing known about lab method

The following rules and best practices should be followed when populating this table:

1. This table will only be populated if data values associated with physical samples are added to the ODM database.
2. The SamplID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
3. The SampleType field should be populated using terms from the SampleTypeCV table. Where the sample type is unknown, a default value of “Unknown” can be used.
4. The LabSampleCode should be a unique text code used by the laboratory to identify the sample. This field is an alternate key for this table and should be unique.
5. The LabMethodID must reference a valid LabMethodID from the LabMethods table. The LabMethods table should be populated with the appropriate laboratory method information prior to adding records to this table that reference that laboratory method. A default value of 0 for this field indicates that nothing is known about the laboratory method used to analyze the sample.

**Table: SampleTypeCV**

The SampleTypeCV table contains the controlled vocabulary for sample type.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for sample type.	“FD”, “PB”, “Grab Sample”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of sample type controlled vocabulary term. The definition is optional if the term is self explanatory.	“Foliage Digestion”, “Precipitation Bulk”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: SeriesCatalog**

The SeriesCatalog table lists each separate data series in the database for the purposes of identifying or displaying what data are available at each site and to speed simple queries without querying the main DataValues table. Unique site/variable combinations are defined by unique combinations of SiteID, VariableID, MethodID, SourceID, and QualityControlLevelID.

This entire table should be programmatically derived and should be updated every time data is added to the database. Constraints on each field in the SeriesCatalog table are dependent upon the constraints on the fields in the table from which those fields originated.

Field Name	Data Type	Description	Example	Constraint
SeriesID	Integer Identity	Unique integer identifier for each data series.	5	P Primary key
SiteID	Integer	Site identifier from the Sites table.	7	P
SiteCode	Text (50)	Site code used by organization that collects the data.	“1002000”	P
SiteName	Text (255)	Full text name of sampling site.	“Logan River”	P
VariableID	Integer	Integer identifier for each Variable that references the Variables table.	4	P
VariableCode	Text (50)	Variable code used by the organization that collects the data.	“00060”	P
VariableName	Text (255)	Name of the variable from the variables table.	“Temperature”	P
Speciation	Text (255)	Code used to identify how the data value is expressed (i.e., total phosphorus concentration expressed <i>as P</i> ). This should be from the SpeciationCV controlled vocabulary table.	“P”, “N”, “NO3”	P
VariableUnitsID	Integer	Integer identifier that references the record in the Units table giving the Units of the data value.	5	P
VariableUnitsName	Text (255)	Full text name of the variable units from the UnitsName field in the Units table.	“milligrams per liter”	P
SampleMedium	Text (255)	The medium of the sample. This should be from the SampleMediumCV controlled vocabulary table.	“Surface Water”	P
ValueType	Text (255)	Text value indicating what type of data value is being recorded. This should be from the ValueTypeCV controlled vocabulary table.	“Field Observation”	P
TimeSupport	Real	Numerical value that indicates the time support (or temporal footprint) of the data values. 0 is used to indicate data values that are instantaneous. Other values indicate the time over which the data values are implicitly or explicitly averaged or aggregated.	0, 24	P

<b>Field Name</b>	<b>Data Type</b>	<b>Description</b>	<b>Example</b>	<b>Constraint</b>
TimeUnitsID	Integer	Integer identifier that references the record in the Units table giving the Units of the time support. If TimeSupport is 0, indicating an instantaneous observation, a unit needs to still be given for completeness, although it is somewhat arbitrary.	4	P
TimeUnitsName	Text (255)	Full text name of the time support units from the UnitsName field in the Units table.	"hours"	P
DataType	Text (255)	Text value that identifies the data as one of several types from the DataTypeCV controlled vocabulary table.	"Continuous" "Instantaneous" "Cumulative" "Incremental" "Average" "Minimum" "Maximum" "Constant Over Interval" "Categorical"	P
GeneralCategory	Text (255)	General category of the variable from the GeneralCategoryCV table.	"Water Quality"	P
MethodID	Integer	Integer identifier that identifies the method used to generate the data values and references the Methods table.	2	P
MethodDescription	Text (unlimited)	Full text description of the method used to generate the data values.	"Specific conductance measured using a Hydrolab" or "Streamflow measured using a V notch weir with dimensions xxx"	P
SourceID	Integer	Integer identifier that identifies the source of the data values and references the Sources table.	5	P
Organization	Text (255)	Text description of the source organization from the Sources table.	"USGS"	P
SourceDescription	Text (unlimited)	Text description of the data source from the Sources table.	"Text file retrieved from the EPA STORET system indicating data originally from Utah Division of Water Quality"	P

<b>Field Name</b>	<b>Data Type</b>	<b>Description</b>	<b>Example</b>	<b>Constraint</b>
Citation	Text (unlimited)	Text string that give the citation to be used when the data from each source are referenced.	“Slaughter, C. W., D. Marks, G. N. Flerchinger, S. S. Van Vactor and M. Burgess, (2001), "Thirty-five years of research data collection at the Reynolds Creek Experimental Watershed, Idaho, United States," Water Resources Research, 37(11): 2819-2823.”	P
QualityControlLevelID	Integer	Integer identifier that indicates the level of quality control that the data values have been subjected to.	0,1,2,3,4	P
QualityControlLevelCode	Text (50)	Code used to identify the level of quality control to which data values have been subjected.	“1”, “1.1”, “Raw”, “QC Checked”	P
BeginDateTime	Date/Time	Date of the first data value in the series. To be programmatically updated if new records are added.	9/4/2003 7:00:00 AM	P
EndDateTime	Date/Time	Date of the last data value in the series. To be programmatically updated if new records are added.	9/4/2005 7:00:00 AM	P
BeginDateTimeUTC	Date/Time	Date of the first data value in the series in UTC. To be programmatically updated if new records are added.	9/4/2003 2:00 PM	P
EndDateTimeUTC	Date/Time	Date of the last data value in the series in UTC. To be programmatically updated if new records are added.	9/4/2003 2:00 PM	P
ValueCount	Integer	The number of data values in the series identified by the combination of the SiteID, VariableID, MethodID, SourceID and QualityControlLevelID fields. To be programmatically updated if new records are added.	50	P

**Table: Sites**

The Sites table provides information giving the spatial location at which data values have been collected.

Field Name	Data Type	Description	Example	Constraint	Default Value
SiteID	Integer Identity	Unique identifier for each sampling location.	37	M Unique Primary key	
SiteCode	Text (50)	Code used by organization that collects the data to identify the site	“10109000” (USGS Gage number)	M Unique Allows only characters in the range of A-Z (case insensitive), 0-9, and “.”, “_”, and “-”.	
SiteName	Text (255)	Full name of the sampling site.	“LOGAN RIVER ABOVE STATE DAM, NEAR LOGAN,UT”	M Cannot contain tab, line feed, or carriage return characters	
Latitude	Real	Latitude in decimal degrees.	45.32	M (>= -90 AND <= 90)	
Longitude	Real	Longitude in decimal degrees. East positive, West negative.	-100.47	M (>= -180 AND <= 360)	
LatLongDatumID	Integer	Identifier that references the Spatial Reference System of the latitude and longitude coordinates in the SpatialReferences table.	1	M Foreign key	0 = Unknown
Elevation_m	Real	Elevation of sampling location (in m). If this is not provided it needs to be obtained programmatically from a DEM based on location information.	1432	O	NULL
VerticalDatum	Text (255)	Vertical datum of the elevation. Controlled Vocabulary from VerticalDatumCV.	“NAVD88”	O Foreign key	NULL
LocalX	Real	Local Projection X coordinate.	456700	O	NULL
LocalY	Real	Local Projection Y Coordinate.	232000	O	NULL

Field Name	Data Type	Description	Example	Constraint	Default Value
LocalProjectionID	Integer	Identifier that references the Spatial Reference System of the local coordinates in the SpatialReferences table. This field is required if local coordinates are given.	7	O Foreign key	NULL
PosAccuracy_m	Real	Value giving the accuracy with which the positional information is specified in meters.	100	O	NULL
State	Text (255)	Name of state in which the monitoring site is located.	“Utah”	O Cannot contain tab, line feed, or carriage return characters	NULL
County	Text (255)	Name of county in which the monitoring site is located.	“Cache”	O Cannot contain tab, line feed, or carriage return characters	NULL
Comments	Text (unlimited)	Comments related to the site.		O	NULL

The following rules and best practices should be followed when populating this table:

1. The SiteID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
2. The SiteCode field must contain a text code that uniquely identifies each site. The values in this field should be unique and can be an alternate key for the table. SiteCodes cannot contain any characters other than A-Z (case insensitive), 0-9, period “.”, dash “-“, and underscore “\_”.
3. The LatLongDatumID must reference a valid SpatialReferenceID from the SpatialReferences controlled vocabulary table. If the datum is unknown, a default value of 0 is used.
4. If the Elevation\_m field is populated with a numeric value, a value must be specified in the VerticalDatum field. The VerticalDatum field can only be populated using terms from the VerticalDatumCV table. If the vertical datum is unknown, a value of “Unknown” is used.
5. If the LocalX and LocalY fields are populated with numeric values, a value must be specified in the LocalProjectionID field. The LocalProjectionID must reference a valid SpatialReferenceID from the SpatialReferences controlled vocabulary table. If the spatial reference system of the local coordinates is unknown, a default value of 0 is used.

**Table: Sources**

The Sources table lists the original sources of the data, providing information sufficient to retrieve and reconstruct the data value from the original data files if necessary.

Field Name	Data Type	Description	Example	Constraint	Default Value
SourceID	Integer Identity	Unique integer identifier that identifies each data source.	5	M Unique Primary key	
Organization	Text (255)	Name of the organization that collected the data. This should be the agency or organization that collected the data, even if it came out of a database consolidated from many sources such as STORET.	“Utah Division of Water Quality”	M Cannot contain tab, line feed, or carriage return characters	
SourceDescription	Text (unlimited)	Full text description of the source of the data.	“Text file retrieved from the EPA STORET system indicating data originally from Utah Division of Water Quality”	M	
SourceLink	Text (500)	Link that can be pointed at the original data file and/or associated metadata stored in the digital library or URL of data source.		O	NULL
ContactName	Text (255)	Name of the contact person for the data source.	“Jane Adams”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
Phone	Text (255)	Phone number for the contact person.	“435-797-0000”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
Email	Text (255)	Email address for the contact person.	“Jane.Adams@dwq.ut”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”



Field Name	Data Type	Description	Example	Constraint	Default Value
Address	Text (255)	Street address for the contact person.	“45 Main Street”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
City	Text (255)	City in which the contact person is located.	“Salt Lake City”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
State	Text (255)	State in which the contact person is located. Use two letter abbreviations for US. For other countries give the full country name.	“UT”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
ZipCode	Text (255)	US Zip Code or country postal code.	“82323”	M Cannot contain tab, line feed, or carriage return characters	“Unknown”
Citation	Text (unlimited)	Text string that give the citation to be used when the data from each source are referenced.	“Data collected by USU as part of the Little Bear River Test Bed Project”	M	“Unknown”
MetadataID	Integer	Integer identifier referencing the record in the ISOMetadata table for this source.	5	M Foreign key	0 = Unknown or uninitialized metadata

The following rules and best practices should be followed when populating this table:

1. The SourceID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
2. The Organization field should contain a text description of the agency or organization that created the data.
3. The SourceDescription field should contain a more detailed description of where the data was actually obtained.
4. A default value of “Unknown” may be used for the source contact information fields in the event that this information is not known.
5. Each source must be associated with a metadata record in the ISOMetadata table. As such, the MetadataID must reference a valid MetadataID from the ISOMetadata table. The ISOMetatadata table should be populated with an appropriate record prior to adding

a source to the Sources table. A default MetadataID of 0 can be used for a source with unknown or uninitialized metadata.

6. Use the Citation field to record the text that you would like others to use when they are referencing your data. Where available, journal citations are encouraged to promote the correct crediting for use of data.

**Table: SpatialReferences**

The SpatialReferences table provides information about the Spatial Reference Systems used for latitude and longitude as well as local coordinate systems in the Sites table. This table is a controlled vocabulary.

Field Name	Data Type	Description	Example	Constraint
SpatialReferenceID	Integer Identity	Unique integer identifier for each Spatial Reference System.	37	M Unique Primary key
SRSID	Integer	Integer identifier for the Spatial Reference System from <a href="http://www.epsg.org/">http://www.epsg.org/</a> .	4269	O
SRSName	Text (255)	Name of the Spatial Reference System.	“NAD83”	M Cannot contain tab, line feed, or carriage return characters
IsGeographic	Boolean	Value that indicates whether the spatial reference system uses geographic coordinates (i.e. latitude and longitude) or not.	“True”, “False”	O
Notes	Text (unlimited)	Descriptive information about the Spatial Reference System. This field would be used to define a non-standard study area specific system if necessary and would contain a description of the local projection information. Where possible, this should refer to a standard projection, in which case latitude and longitude can be determined from local projection information. If the local grid system is non-standard then latitude and longitude need to be included too.		O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: SpeciationCV**

The SpeciationCV table contains the controlled vocabulary for the Speciation field in the Variables table.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for Speciation.	“P”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of Speciation controlled vocabulary term. The definition is optional if the term is self explanatory.	“Expressed as phosphorus”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: TopicCategoryCV**

The TopicCategoryCV table contains the controlled vocabulary for the ISOMetaData topic categories.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for TopicCategory.	“InlandWaters”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of TopicCategory controlled vocabulary term. The definition is optional if the term is self explanatory.	“Data associated with inland waters”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: Units**

The Units table gives the Units and UnitsType associated with variables, time support, and offsets. This is a controlled vocabulary table.

Field Name	Data Type	Description	Example	Constraint
UnitsID	Integer Identity	Unique integer identifier that identifies each unit.	6	M Unique Primary key
UnitsName	Text (255)	Full text name of the units.	“Milligrams Per Liter”	M Cannot contain tab, line feed, or carriage return characters
UnitsType	Text (255)	Text value that specifies the dimensions of the units.	“Length” “Time” “Mass”	M Cannot contain tab, line feed, or carriage return characters
UnitsAbbreviation	Text (255)	Text abbreviation for the units.	“mg/L”	M Cannot contain tab, line feed, or carriage return characters

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: ValueTypeCV**

The ValueTypeCV table contains the controlled vocabulary for the ValueType field in the Variables and SeriesCatalog tables.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for ValueType.	“Field Observation”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of the ValueType controlled vocabulary term. The definition is optional if the term is self explanatory.	“Observation of a variable using a field instrument”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: VariableNameCV**

The VariableName CV table contains the controlled vocabulary for the VariableName field in the Variables and SeriesCatalog tables.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for Variable names.	"Temperature", "Discharge", "Precipitation"	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of the VariableName controlled vocabulary term. The definition is optional if the term is self explanatory.		O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.

**Table: Variables**

The Variables table lists the full descriptive information about what variables have been measured.

Field Name	Data Type	Description	Example	Constraint	Default Value
VariableID	Integer Identity	Unique integer identifier for each variable.	6	M Unique Primary key	
VariableCode	Text (50)	Text code used by the organization that collects the data to identify the variable.	"00060" used by USGS for discharge	M Unique Allows only characters in the range of A-Z (case insensitive), 0-9, and ".", "-", and "_".	
VariableName	Text (255)	Full text name of the variable that was measured, observed, modeled, etc. This should be from the VariableNameCV controlled vocabulary table.	"Discharge"	M Foreign key	

Field Name	Data Type	Description	Example	Constraint	Default Value
Speciation	Text (255)	Text code used to identify how the data value is expressed (i.e., total phosphorus concentration expressed as <i>P</i> ). This should be from the SpeciationCV controlled vocabulary table.	“P”, “N”, “NO3”	M Foreign key	“Not Applicable”
VariableUnitsID	Integer	Integer identifier that references the record in the Units table giving the units of the data values associated with the variable.	4	M Foreign key	
SampleMedium	Text (255)	The medium in which the sample or observation was taken or made. This should be from the SampleMediumCV controlled vocabulary table.	“Surface Water” “Sediment” “Fish Tissue”	M Foreign key	“Unknown”
ValueType	Text (255)	Text value indicating what type of data value is being recorded. This should be from the ValueTypeCV controlled vocabulary table.	“Field Observation” “Laboratory Observation” “Model Simulation Results”	M Foreign key	“Unknown”
IsRegular	Boolean	Value that indicates whether the data values are from a regularly sampled time series.	“True” “False”	M	“False”
TimeSupport	Real	Numerical value that indicates the time support (or temporal footprint) of the data values. 0 is used to indicate data values that are instantaneous. Other values indicate the time over which the data values are implicitly or explicitly averaged or aggregated.	0, 24	M	0 = Assumes instantaneous samples where no other information is available
TimeUnitsID	Integer	Integer identifier that references the record in the Units table giving the Units of the time support. If TimeSupport is 0, indicating an instantaneous observation, a unit needs to still be given for completeness, although it is somewhat arbitrary.	4	M Foreign key	103 = hours
DataType	Text (255)	Text value that identifies the data values as one of several types from the DataTypeCV controlled vocabulary table.	“Continuous” “Sporadic” “Cumulative” “Incremental” “Average” “Minimum” “Maximum” “Constant Over Interval” “Categorical”	M Foreign key	“Unknown”

Field Name	Data Type	Description	Example	Constraint	Default Value
GeneralCategory	Text (255)	General category of the data values from the GeneralCategoryCV controlled vocabulary table.	“Climate” “Water Quality” “Groundwater Quality”	M Foreign key	“Unknown”
NoDataValue	Real	Numeric value used to encode no data values for this variable.	-9999	M	-9999

The following rules and best practices should be followed when populating this table:

1. The VariableID field is the primary key, must be a unique integer, and cannot be NULL. This field should be implemented as an auto number/identity field.
2. The VariableCode field must be unique and serves as an alternate key for this table. Variable codes can be arbitrary, or they can use an organized system. VariableCodes cannot contain any characters other than A-Z (case insensitive), 0-9, period “.”, dash “-“, and underscore “\_”.
3. The VariableName field must reference a valid Term from the VariableNameCV controlled vocabulary table.
4. The Speciation field must reference a valid Term from the SpeciationCV controlled vocabulary table. A default value of “Not Applicable” is used where speciation does not apply. If the speciation is unknown, a value of “Unknown” can be used.
5. The VariableUnitsID field must reference a valid UnitsID from the UnitsTable controlled vocabulary table.
6. Only terms from the SampleMediumCV table can be used to populate the SampleMedium field. A default value of “Unknown” is used where the sample medium is unknown.
7. Only terms from the ValueTypeCV table can be used to populate the ValueType field. A default value of “Unknown” is used where the value type is unknown.
8. The default for the TimeSupport field is 0. This corresponds to instantaneous values. If the TimeSupport field is set to a value other than 0, an appropriate TimeUnitsID must be specified. The TimeUnitsID field can only reference valid UnitsID values from the Units controlled vocabulary table. If the TimeSupport field is set to 0, any time units can be used (i.e., seconds, minutes, hours, etc.), however a default value of 103 has been used, which corresponds with hours.
9. Only terms from the DataTypeCV table can be used to populated the DataType field. A default value of “Unknown” can be used where the data type is unknown.
10. Only terms from the GeneralCategoryCV table can be used to populate the GeneralCategory field. A default value of “Unknown” can be used where the general category is unknown.
11. The NoDataValue should be set such that it will never conflict with a real observation value. For example a NoDataValue of -9999 is valid for water temperature because we would never expect to measure a water temperature of -9999. The default value for this field is -9999.

**Table: VerticalDatumCV**

The VerticalDatumCV table contains the controlled vocabulary for the VerticalDatum field in the Sites table.

Field Name	Data Type	Description	Examples	Constraint
Term	Text (255)	Controlled vocabulary for VerticalDatum.	“NAVD88”	M Unique Primary key Cannot contain tab, line feed, or carriage return characters
Definition	Text (unlimited)	Definition of the VerticalDatum controlled vocabulary. The definition is optional if the term is self explanatory.	“North American Vertical Datum of 1988”	O

This table is pre-populated within the ODM. Changes to this controlled vocabulary can be requested at <http://water.usu.edu/cuahsi/odm/>.



## Appendix B. Data Versioning Within ODM

The main text of this document focuses on how ODM is structured to store observations data. It does not address how to manage editing data stored within ODM. Software applications based on ODM will have functionality that will allow data managers and database administrators to modify, delete, change, or otherwise make edits to data stored within ODM. In addition, these software tools will provide functionality to create derived datasets, or datasets that are calculated or derived from data already stored in ODM (i.e., calculate a time series of discharge from a time series of stage, or calculate a time series of daily average temperature from a time series of hourly observations). The purpose of this appendix is to clarify how data editing and versioning can be managed within the ODM schema.

### Data Series Defined

In order to fully grasp the concepts that follow, the idea of a “data series” in the context of ODM must be clarified. A “data series” is an organizing principle that is present in the ODM. A data series consists of all of the data values associated with a unique site, variable, method, source, and quality control level combination. An example of the full specification for a data series is: “all of the raw unchecked (QualityControlLevel) water temperature (Variable) values measured in the Logan River near Logan, UT (Site) using a field temperature sensor (method) by Utah State University (Source).” Each record in the SeriesCatalog table of ODM represents a unique data series.

### Rules for Editing and Deriving Data Series in ODM

The following rules are suggested so that versioning of and edits to data series can be managed within the ODM schema. Software applications that work with ODM should follow these rules. These rules are based on the default set of Quality Control Levels that are distributed with the ODM blank schema.

1. *Data versioning should be done at the data series level* – Within ODM, the concept of data versioning is related to the quality control level. Quality control level is a data series level attribute, and as such, changes to the quality control level should occur at the data series level rather than at the individual value level. For example, if an investigator wished to create a quality controlled Level 1 data series from a raw Level 0 data series, he/she should first make a copy of the raw Level 0 data series and then perform any edits and adjustments required in the quality control process to the copy. The edited copy then becomes the Level 1 data series, and the Level 0 data series is preserved intact.
2. *Data series with a QualityControlLevelCode of 0 cannot be edited* – Level 0 data series represent raw data from sensors (i.e., stage measured by a water level recorder) or other products derived from raw data (i.e., discharge that is programmatically derived from stage before the stage values have been quality controlled). By definition, Level 0 data have not been quality controlled and may contain significant errors and bad values. However, Level 0 data series represent the source from which all other derived data series are based, and as such should remain intact for archive purposes. Level 0 data series should not be used for analysis unless no other adequate options are available, and

only if the user is aware that the data are raw. Level 0 data series can be removed entirely from the database, but only by removing the entire data series.

3. *Only one QualityControlLevel 0 data series can exist for a Site, Variable, and Method combination* – Only one raw data series for a Site, Variable, and Method combination can exist within an ODM database. If multiple sensors are measuring the same variable at the same site, the method description would have to distinguish between the two.
4. *Only one QualityControlLevel 1 data series can exist for each Site, Variable, and Method combination* – Once a Level 0 data series has been loaded to the database, a Level 1 data series can be “derived” from that Level 0 data series. This is done by first making a copy of the Level 0 data series, second changing the QualityControlLevel of the copy to 1, and last doing any necessary filtering or editing required so that the Level 1 data series is acceptable as quality controlled. In most cases, the majority of the values within a Level 0 data series and its corresponding Level 1 data series will remain the same. However, where instruments malfunction or other conditions are present that affect the raw data values, Level 0 values may be deleted, adjusted, or otherwise edited in creating the Level 1 data series.
5. *Any edits to a data series are saved to that data series* – Level 0 data cannot be edited. With Levels 1 or higher, however, software applications should be allowed to edit and delete values. Each time an edit is made, the result should overwrite the previous value within a data series. In other words, edits should not create new data series, they should modify an existing one. This will be true even where edits are done within multiple editing sessions. The editing software should record the method or basis for any data edits in appropriate method records.
6. *Data series of Level 2 or higher can only be created from data series of Level 1 or higher* – Derived data series of Level 2 or higher can only be created from data series of Level 1 or higher. If a user wishes to create a derived data series from a Level 0 data series (such as discharge from raw, unchecked stage values) that derived data series would also be Level 0.