

# Hydrologic data management using RDBMS with differential-linear data storage

G. Tsakalias, D. Koutsoyiannis

*Division of Water Resources, Hydraulics and Maritime Engineering, National Technical University of Athens, 5 Iroon Polytechniou, GR-157 00 Zografou, Greece*

## Abstract

Recently, Relational Data Base Management Systems (RDBMSs) have become popular for handling hydrologic time-series and running hydrologic applications. However, the standard data independence in such systems has many technical disadvantages in storing time-series data. An alternative technique, named Differential-Linear Data Storage (DLDS) technique, has been developed in the framework of the Hydroscope project (a Greek nation-wide database for hydrometeorological information). This technique establishes a standardised representation of hydrologic time-series in a relational database environment, also providing a notable reduction in storage space. Instead of the standard SQL queries (select, insert, delete and update), numerous composite procedures are implemented to facilitate time-series management.

## 1. Introduction

During the last three decades, hydrologic data storage and retrieval has been continuously improved as a result of the evolution of computer technologies in data management. The first systems developed were based on central mainframe computers and programming languages like FORTRAN. Recently, new hydrologic database systems have been developed taking advantage on the available new technologies such as Relational Data Base Management Systems (RDBMSs), graphical user interfaces, etc. Examples of such systems are the National Water Information System II in USA (USGS<sup>1</sup>), CompuMod in Canada (Environment Canada<sup>2</sup>) and Hydroscope in Greece (Tolikas et al.<sup>4</sup>). The latter is a nation-wide distributed database system for hydrometeorological information currently consisting of 12 nodes, each belonging to one of the co-operating or-

ganisations (ministries, universities and research centres). It is based on Hewlett-Packard 9000/s700 UNIX workstations running the INGRES RDBMS.

In a RDBMS the data are organised in tables. Each table may be viewed as a rectangular entity consisting of rows (or records) and columns (or fields). Each row is a separate instance of the same kind of information, while the different fields of the row contain different components of this information. One major feature of RDBMSs is data independence: the system itself controls the details of data storage and retrieval, while the user need not know how the data are stored. Another important feature is that they support the use of special programming languages for data manipulation such as SQL (Structured Query Language), which provides powerful procedures for data access and retrieval (i.e., select, insert, delete and update). Owing to these features, RDBMSs have dominated in several business applications and they have become the state-of-the-art in commercially available database technologies. This has further enhanced their capacity, as they are continually improved offering extensive data handling capabilities and support and development tools. Thus, they currently support distributed databases and provide easy-to-use graphical user interfaces and fourth generation programming languages (such as Windows 4GL) and object oriented techniques (Papakostas et al.<sup>3</sup>)

However, RDBMSs have certain disadvantages when used for hydrologic data management (Dodson<sup>1</sup>). These are associated with the time-series character of hydrologic data, their numeric aspect and the complicated calculations needed for certain hydrologic applications. The data independence results in a significant increase of the storage volume, as the date of each data value has to be stored individually with the data value (herein the term "date" means the time of the measurement with the precision offered by the particular RDBMS). In addition, the typical calculations needed for manipulating hydrologic time-series (e.g., interpolation, integration, extraction of maxima or minima) are mainly procedural calculations and thus are not directly supported by non-procedural languages such as SQL. Finally, as it will be described later, the consecutive data values of a time-series are in fact connected to each other such that a possible update of one value may also affect its previous and next data values. This is inconsistent with the data-independence feature of a RDBMS.

Despite these disadvantages, the above mentioned strong features of the commercialised RDBMSs can justify their choice for hydrologic applications. This is the case for the Hydroscope project where the selection of the relational model offered a lot of powerful characteristics that lead to an efficient distributed hydrologic database system. To overcome some of the disadvantages of the relational model, a specific technique, named Differential-Linear Data Storage (DLDS) technique, has been developed and incorporated into the database. The objectives of this technique were (a) to standardise the representation of a time-series in a RDBMS environment; (b) to implement a generalised scheme of a time-series with a varying time step, taking advantage of the data-independence feature of the RDBMS; (c) to cope with the intermittent nature of time-series; (d) to create a set of SQL composite queries that substitute the

typical simple queries in a manner that they are consistent with the data interdependence within a time-series; (e) to provide a set of procedures for typical calculations of a time-series; and (f) to provide an efficient compression technique for time-series data. It is noted that in the Hydroscope system the latter objective (data compression) is also served by other techniques such as integer storage of floating-point values (this is also used in other systems; see Dodson<sup>1</sup>), substitution of the typical date data type offered by the RDBMS with a user-defined data type of a smaller length, and the use of an alternative storage technique for time-series with constant time step, the list storage structure. A brief description of those alternative techniques may be found elsewhere (Papakostas et al.<sup>3</sup>).

## 2. Classification of time-series

Depending on the physical meaning and the observation technique of each hydrologic time-series we can classify hydrologic variables in three categories:

*Instantaneous variables:* Most hydrologic variables such as river stage, river discharge, and temperature are measured (or estimated) and interpreted at a specific time instant (Figure 1a).

*Cumulative variables:* By continuously integrating an instantaneous variable with respect to time, starting at a certain time instant (Figure 1b) we get another type of variable, the cumulative type. A variable of this type is the water volume in a lake or reservoir (the corresponding instantaneous variable being the net inflow into the lake).

*Differential variables:* By intermittently integrating an instantaneous variable with respect to time, starting at numerous arbitrary time instants (Figure 1c) we get a third type of variable, the differential type. Time-series of this type are those of the rainfall depth (each measurement of the rainfall depth is the integral of the intensity for a certain period), the evaporation depth and the sunshine duration.

We can observe that in the first two types the value of the time-series  $x(t)$  for any time instant  $t$  is unique and has a specific physical meaning. However, in the case of differential variables this is not true, because  $x(t)$  depends on the selection of the time instants of discontinuity (i.e., points where the renewal of integration takes place). Consequently, from the point of view of the database manipulation, the differential type needs a different treatment while the other two variable types, which hereinafter will be referred to as non-differential types, can be treated with a common manner. The only difference between instantaneous and cumulative variables is that in the former case given a time period  $(t - \Delta t, t)$  the time average

$$x_{\Delta t}(t) = \frac{1}{\Delta t} \int_{t-\Delta t}^t x(t) dt \quad (1)$$

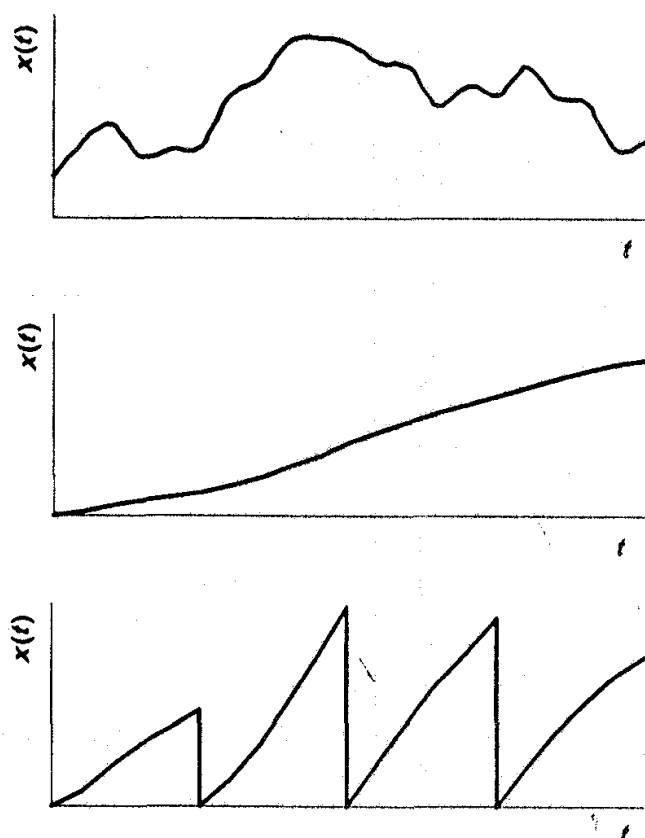
is meaningful whereas in the latter case it is not. A meaningful quantity for that time period in the case of cumulative variables is the difference

$$x'_{\Delta t}(t) = x(t) - x(t - \Delta t) \quad (2)$$

The corresponding quantity in the case of differential variables is

$$x'_{\Delta t}(t) = x(t) - x(t - \Delta t) + \sum_{t - \Delta t < t_i < t} x(t_i) \quad (3)$$

where  $t_i$  is any time instant of discontinuity (renewal of integration).



**Figure 1.** Explanation sketch for the classification of time-series: (a) instantaneous type, (b) cumulative type, and (c) differential type.

Because a time-series is never known in continuous time, we will assume that  $x(t)$  is known at a set of points  $t_j$  only, not necessarily equidistant. Also, we will assume that the distances between consecutive points are the largest distances that can guarantee a valid linear interpolation between them. The series  $x(t_j)$  is stored in the database and can be interpreted in two ways, depending on the type of variable (differential or non-differential) as indicated in Figure 2. The technique for the storage and interpretation of  $x(t)$  that is described in the next sections, will be referred to as the Differential-Linear Data Storage (DLDS) technique, because of the assumption of linear interpolation and its ability to handle differential time-series.

### 3. Data interdependence

We can observe from Figures 1 and 2 that there is always an interdependence between the separate rows of a time-series in a database table. This is especially apparent in the differential variables where the value at date  $t$  has meaning only if we know the date of the previous measurement. Similarly in the case of non-differential time-series we should know when we are not allowed to interpolate any more (when there is a lack of information). For example we see in Figure 2a that we can interpolate (or perform some other calculation) between  $t_2$  and  $t_3$ , but not between  $t_3$  and  $t_5$  because there is lack of information in that period (e.g., due to instrument malfunction). Thus, it is necessary to add a record with a null value between  $t_3$  and  $t_5$  (a null is a special value indicating the absence of a valid observation), which will indicate that interpolation is not allowed.

It is obvious that in situations like the described above, simple tasks like the deletion or insertion of a measurement cannot be performed without the consideration of the data interdependence. We cannot simply delete a measurement, but it may be necessary to replace the deleted record with a null value record to indicate that an interpolation or integration at the measurement's neighbourhood is no longer valid. We also cannot simply insert a value because this will probably alter other values computed by interpolation or integration.

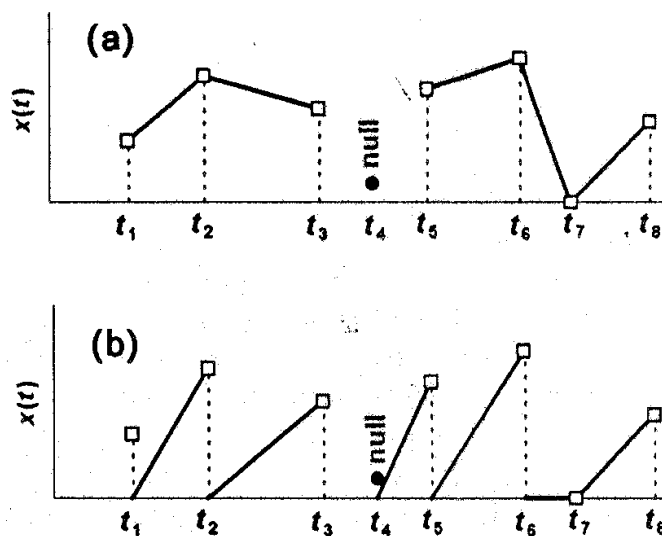


Figure 2 Explanation sketch for the interpretation of stored data series: (a) non-differential series and (b) differential series

### 4. DLDS chain definition

In relational databases the records are not ordered, thus the notion of the *next* row does not exist. We will explicitly define here the *next* record of the record with date  $d_n$  as the record with date  $d_{n+1}$  that can be found as:

```
select  $d_{n+1} = \min(\text{date})$  where date  $\geq d_n$ 
```

In the above notation (as well as in following similar situations) we use an SQL like pseudocode to enhance readability. In a similar way we can define the *previous* date as the date  $d_{n-1}$ . We will refer to a record with date  $d_n$  as the *record<sub>n</sub>*.

A DLDS chain is a set of records in a database table representing a part of a time-series and having at least two columns standing for the date and the value of an observation. If record<sub>1</sub> is the member of the chain with the minimum date and record<sub>n</sub> is the record with the maximum date then: (a) the record<sub>1</sub> is the only member of the chain that has a null value; (b) if  $d_1 < d_m < d_n$  then record<sub>m</sub> is a member of the chain; and (c) the set of records of the chain is not subset of any set of records that form another DLDS chain.

The chain can be interpreted as a differential or a non-differential one (see Figure 2). Obviously a time-series consists of one or more DLDS chains (the number of chains depends on the number of periods with lack of information)

## 5. User-programme interaction

We will examine now how a computer user can manipulate a DLDS chain. Obviously he or she cannot use the conventional procedures used in RDBMSs due to the implicit data interdependence. Despite this, we should consider that the user is familiar with a standard way of interaction with a database front end. The question is if we can keep this type of interaction and use it to manipulate DLDS chains.

We will examine the typical situation where we have a tablefield (an array of fields on the computer screen) that displays time-series data and a user performing actions on the tablefield. We will also assume that the database table (from which the data come from) has two columns only, date and value, and that the column date is a unique key. The tablefield has three columns, date, value and hidden-date. The user can change the value of columns date and value while he or she cannot see the third column.

Each row of the tablefield has a row state (or a destination) associated with it. Every user action modifies this destination. At the end of the session when a *save* command is executed the row destinations will lead to specific database queries. In an independent row format the modification of the row destinations according to user actions and the interpretation of the row destinations to database queries could be briefly described as in Table 1.

In all cases before the insert command the save procedure checks if a record with the same date already exists and provides the user a dialog box to decide if the existing record should be overwritten.

We can observe in Table 1 that a row destination always leads to a unique SQL query. Hence, we want to keep the user-programme interactions unchanged (thus retaining the consistency of the user interface) and use these row destinations to manipulate DLDS chains. This can be done by substituting each single query (update, insert, delete) with a sequence of queries, depending on the chain surrounding the processed row, as we can see in the next section.

Table 1. User-programme interaction

Previous row state	User action	Tablefield modification	New row state	Database query
none	select from the database	the rows appear on the tablefield	unchanged	none
unchanged or changed	write new value		changed	update
unchanged or changed	write new date	an invisible row inheriting the hidden date of the parent row is added	new deleted	insert delete
unchanged or changed	delete row	the row becomes invisible	deleted	delete
none	insert new row		undefined <sup>f</sup>	none
undefined or new	write new value		new	insert
undefined or new	write new date		new	insert
new	delete row	the row becomes invisible	deleted	delete

## 6. Mapping of standard SQL to DLDS queries

The mapping of the single queries (update, insert, delete) to composite DLDS queries is given in Table 2. Each composite query is as a set of rules, all triggered each time a row is considered by the save procedure. Each rule has a condition part and an action part. Generally, the rules are not mutually exclusive. The composite queries keep the name of the replaced single queries with the addition of the character '#'. Some columns of Table 2 need explanation:

*Chain conditions:* In order for a rule to perform its action some conditions must hold. The chain conditions refer to the processed row and the surrounding rows of the DLDS chain.

*Rule conditions:* Some rules are mutually exclusive. This is explicitly declared in this column.

*Dialog:* Before a rule's action is performed, a dialog box may appear. This dialog box can be one of two types: *Warning* is just a message providing the user useful information. *Confirmation* provides the user a message and he or she can answer 'yes' or 'no'. If the user answers 'yes' the rule's action will be performed otherwise no action will be taken.

We note here that in the case of the first insertion in the database the rule I6 creates the first DLDS chain. We also note that the actions taken after the update# DLDS query are also DLDS queries.

Table 2. DLDS queries that replace standard SQL queries

DLDS Query	Rule Name	Chain Conditions	Rule Conditions	Dialog	Action
delete# the record <sub>n</sub>	D1	value <sub>n</sub> is not null			update value <sub>n</sub> to null
	D2	value <sub>n</sub> is null value <sub>n+1</sub> is null			delete record <sub>n</sub>
	D3	value <sub>n</sub> is null value <sub>n+1</sub> is not null		confirmation "The chain will be destroyed. You should probably insert a record to maintain the chain. Should the record be deleted?" [yes]	delete record <sub>n</sub>
insert# a record with date $d_n$	I1	exists record <sub>n</sub> with date $d_n$ ( $d_n$ is a junction of two links at the chain)		confirmation "Should the record <sub>n</sub> be overwritten?" [yes]	delete existing record <sub>n</sub> insert new record <sub>n</sub>
	I2	value <sub>n+1</sub> is null	not I1 not I4		insert record <sub>n</sub>
	I3	exists record <sub>a</sub> , $d_a < d_n$ exists record <sub>b</sub> , $d_b > d_n$ ( $d_n$ is in a chain range and lies in the middle of a link)	not I1	confirmation "The link will be split into the links ( $d_a-d_n$ , $d_n-d_b$ ) thus modifying the chain. Should the new record be inserted?" [yes]	insert record <sub>n</sub>
	I4	not exists record <sub>b</sub> , $d_b > d_n$ (the new record lies after the last chain)	not I1		insert record <sub>n</sub>
	I5	value <sub>n</sub> is null not exists record <sub>a</sub> , $d_a < d_n$ (it is the "first" record in the table)	not I1		insert record <sub>n</sub>
	I6	value <sub>n</sub> is not null not exists record <sub>a</sub> , $d_a < d_n$ (it is the "first" record in the table)	not I1	warning "The record <sub>n-u</sub> with a null value will be inserted just before <sup>§</sup> this record <sub>n</sub> to indicate the start of the chain."	insert record <sub>n</sub> insert record <sub>n-u</sub>
update# a record with date $d_n$ and hidden date $d_{nh}$	U1				delete# record <sub>nh</sub> insert# record <sub>n</sub>

<sup>§</sup> Since time is discrete in DBMSs "just before" means one time unit (u) before



The standard select query (not contained in Table 2) is also substituted by sequences of queries. As there are many forms of it, we will give here one example only. The procedure

```
select X = value where date = d and range = r
```

for a differential variable will be substituted by the following SQL queries (in pseudocode):

```
select da = max(date) where date <= d - r
select db = min(date) where date >= d - r
select valuea = value where date = da
select valueb = value where date = db
value1 = interpolation(valuea, valueb, da, d, db)
select dc = max(date) where date <= d
select dd = min(date) where date >= d
select valuec = value where date = dc
select valued = value where date = dd
value2 = interpolation(valuec, valued, dc, d, dd)
select i = sum(value) where date > d-r and date < d
X = i - value1 + value2
```

We note that if value<sub>1</sub> or value<sub>2</sub> or i is null then the result X will be also null. By a way similar to the procedure above we can define a select procedure for mean non-differential variables, etc.

## 7. Concluding remarks

In summary, the DLDS technique standardises the representation of a time-series in a RDBMS environment, implements a generalised scheme of a time-series with a varying time step, and provides a set of composite SQL queries, which substitute the standard ones. In addition to the above, we have developed a number of procedures that facilitate typical time-series calculations such as transformation of a fixed step time-series to a DLDS chain (and vice versa), extraction of time averaged series, extraction of extreme values, etc.

Due to its structural features (variable time step, linear variation assumption) the DLDS technique results in an efficient compression of time-series data. To give an indication of the reduction in storage space we consider the case of a daily rainfall time-series. Assuming that we have an average of 50 rainy days per year with about 30 sequences of dry days (several days each) between them (a situation typical in Greece), with the DLDS technique we must store  $50 + 30 = 80$  daily values per year (instead of 365). This results in a reduction in space equal to  $100 * (1 - 80/365) = 78\%$ . As another interesting case we consider a time-series of hourly river discharge. Assuming that we have about 20 days per year with significant flood events, while in the other days the variation of discharge is small, such that it suffices to interpolate between daily values to find the hourly values, we need  $20 * 24 + (265 - 20) = 825$  data values in total. In case of the common storage method the total values would be  $365 * 24 = 8760$ . Thus the technique leads to a  $100 * (1 - 825/8760) \approx 90\%$  reduction in storage.

However, the reduction is not so impressive for other variables such as temperature, humidity and so.

The DLDS technique is more effective in data compression than another similar method, the repeat counter technique (consisting of the use of a counter in each row indicating how many repeated values in sequence exist in a time-series; see Dodson<sup>1</sup>). The repeat counter field leads to an about 40% increase of the storage volume for each row. Thus for the daily rainfall example, with this technique the reduction in storage would be  $100 * (1 - 1.40 * 80/365) = 69\%$ . Moreover, the latter technique does not take advantage of linearly varying time-series, such as the river discharge during days with no flood. Finally, the DLDS technique provides a more natural way in data entry and retrieval, as the use of a varying time step is more favourable for the user when compared to the use of a counter, especially when digitising charts (e.g. rain recorder charts) to enter them to the database.

*Acknowledgement:* The authors wish to thank N. Papakostas for his constructive comments.

## References

1. Dodson, R. D., Advances in hydrologic computation, Chapter 23, *Handbook of hydrology*, ed D. R. Maidment, McGraw-Hill, 1992.
2. Environment Canada, *CompuMod functional specification*, Report, Environment Canada, Ottawa, Canada, 1991.
3. Papakostas, N., Nalbantis, I. & Koutsoyiannis, D., Modern computer technologies in hydrologic data management, *Proceedings of the 2<sup>nd</sup> European Conference in Water Resources Technology and Management*, Lisbon, Portugal, 1994.
4. Tolikas, D., Koutsoyiannis, D. & Xanthopoulos, Th., HYDROSCOPE: Un système d'informations pour l'étude des phénomènes hydroclimatiques en Grèce, *6<sup>ème</sup> Colloque International de Climatologie*, Thessaloniki, Grèce, 1993.
5. United States Geological Survey (USGS), *System requirements specification for the USGS National Water Information System II*, Report, S. B. Mathey (ed), USGS, Reston, USA, 1991.