

# **Data Management in the 21<sup>st</sup> Century Emerging Technologies and their Implication for Hydrography**

**Robert M. STIRLING, United Kingdom**

**Key words:** Hydrography, Data, Standards, ISO 15926, EPISTLE, ISPDM, XML, GML.

## **ABSTRACT**

This paper will present one possible 21<sup>st</sup> century approach to the management and exchange of hydrographic data. This involves the use of new and emerging technologies including ISO SC4 data standards for data integration, XML standards for data transport and transformation, data warehousing, distributed system architectures such as J2EE and the adoption of some fundamental data management principles.

The ISPDM project will be used to discuss and outline the implementation of such technologies. ISPDM is a “Take-up” action, funded under the EC’s IST programme, to trial emerging data standards, warehousing and web technologies in the pipeline domain. The project commenced on 1 January 2001, will run eighteen months and is receiving 1,000,000 Euro from the Commission. The project consortium comprises of Thales GeoSolutions (UK), PrismTech (UK), Andrew Palmer & associates (UK), POSC Caesar Association (Norway) and Rosen Engineering (Germany).

The primary benefits from such a data management approach are increased organisational efficiency and effectiveness leading to enhanced decision making and the opportunity to develop new business processes. This results in cost reduction and improved competitiveness. For the hydrographic community the major advantages are:-

- Significant cost savings.
- Management of life-cycle data.
- Vendor independence vis-à-vis application suppliers.
- Resource maximisation.
- An enhanced and more competitive vendor base.
- System future proof - ISO standards are stable and technology independent.
- Data model supported by ISO.
- Web enabled systems.

Maritime safety will also be enhanced through the ready exchange of reliable data between HO’s and between HO’s and ENC vendors, in a consistent manner.

## **CONTACT**

R. M. (Bob) Stirling  
ISPDM Project Manger  
Thales Geosolutions Group Ltd  
12 Beaulieu Crescent (private)  
Kilmacolm  
Scotland PA13 4LR  
UNITED KINGDOM  
Tel. + 44 1505 873 563  
Fax + 44 1505 873 563  
E-mail: BStirlingISPDM@cs.com

# Data Management in the 21<sup>st</sup> Century Emerging Technologies and their Implication for Hydrography

Robert M. STIRLING, United Kingdom

## 1. PRINCIPLES FOR DATA MANAGEMENT AND EXCHANGE IN THE 21<sup>ST</sup> CENTURY

Many current data management systems and exchange formats have their origins in 1970's or 1980's technology. Although many have been revamped to have a "Windows" look and feel the underlying technologies are past their "sell by date". Their continued use is due to vested interests and the limited appreciation of some of the emerging technologies that are being applied to data management.

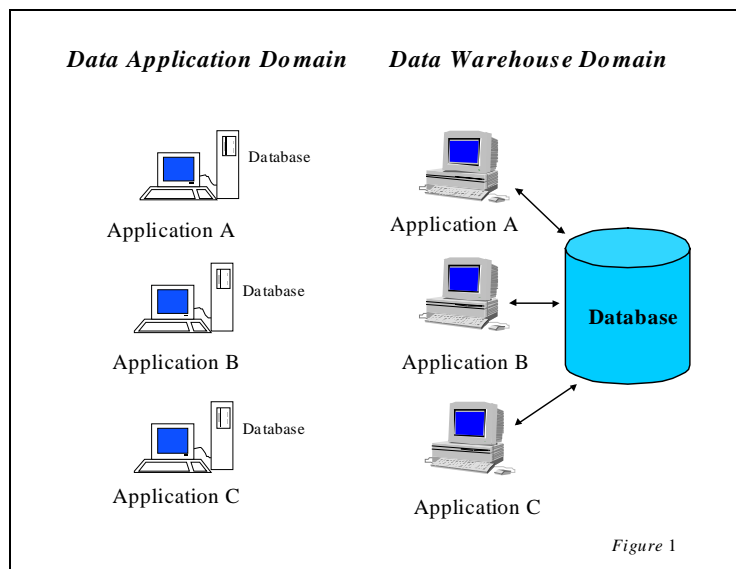
### 1.1 Data Management Domains

There are two basic data management domains:-

*The Data Application Domain.*

*The Data Warehouse Domain.*

In the *application domain* the database and data management functions are an integral part of the application. In the *warehouse domain* the database and data management function is independent of any application. Figure 1 summarises these two domains.



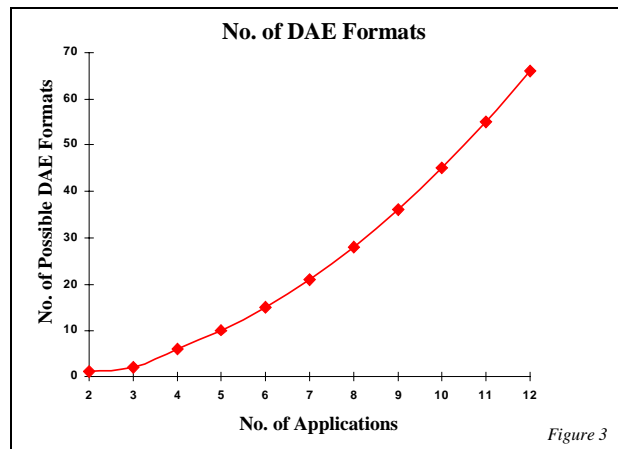
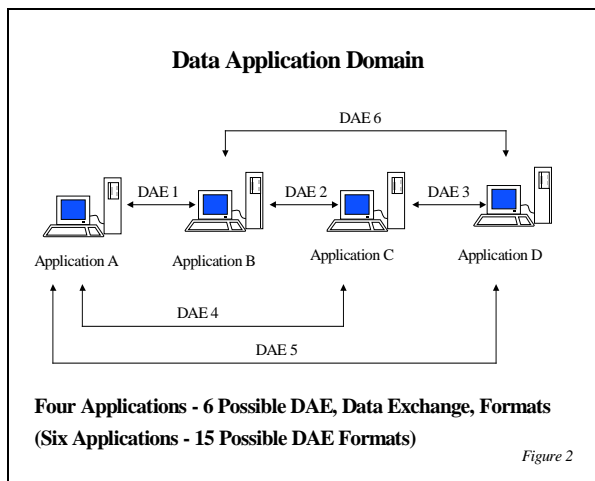
The data application domain has been the most common form for many years. Its prominence has been due to both technical and commercial reasons:-

- Through the 1970's and 1980's the cost and complexity of hardware technology, mainframes, precluded all but the largest organisations establishing large databases.
- As PC power increased it became possible to manage increasingly more complex databases at reasonable cost.
- Limitations of both hardware and software technology meant that the database's data structures had to be designed to optimise the performance of a specific application, or suite of related applications.
- A bespoke database, within an application, often created a competitive advantage for a vendor.

These latter two points create a situation whereby data is "locked in" to an application and cannot be accessed by other applications. This may also lock the client to the vendor, a situation many clients now find unacceptable.

## 1.2 Data Exchange Issues

A particular facet of past, and current, data management is that of data exchange due, primarily, to the fact that the underlying databases had bespoke data models. Exchange of data, between applications, can become quite complex. Data exchange formats, DAE's, have to be developed and continually updated as the applications and supporting databases are modified and upgraded. Figure 2 indicates the type of problem while figure 3 indicates the potential scale of the problem.



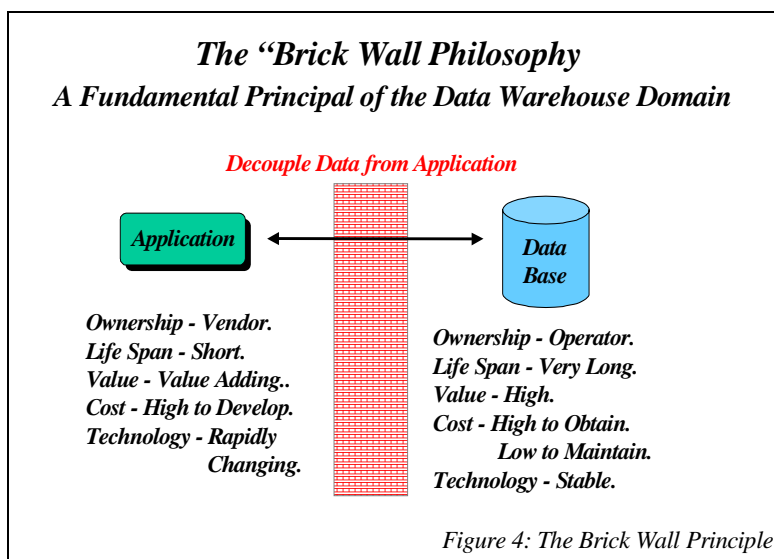
From figure 3 it can be seen that if ten organisations, each with a different database application, wish to exchange data then potentially 45 different data exchange formats are required. This problem has been partly addressed by the development of standards such as S57 and the various SEG and UKOOA formats. These are solutions driven by the technology of the day. They are often limited in content and cannot accommodate the wealth of data required for the modern risk, integrity and asset management applications demanded by a

wide range of industries.

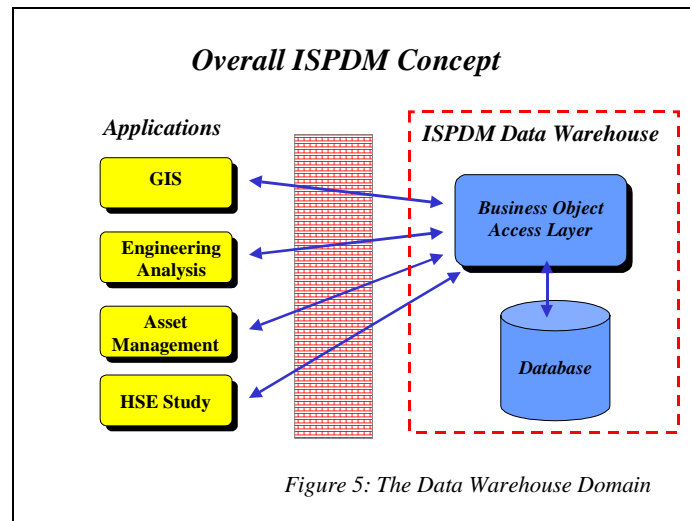
The solution, to the inherent failings of *application domain* data management, is to manage information independent of any application so that it is available to a wide variety of users and applications. This approach is the *data warehouse domain* data management environment.

### 1.3 Data Warehouse Domain

The solution, to the inherent failings of *application domain* data management, is to manage information independent of any application so that it is available to a wide variety of users and applications. This approach is the *data warehouse domain* data management environment. The “Brick Wall” principle, figure 4, is fundamental to a data warehouse domain. This principal is essential as applications and data, hence data management systems, are diametrically different in terms of ownership, life span, value, technology and the skills and knowledge base of the people responsible for each.



To achieve information sharing in the data warehouse domain it is necessary to establish and agree common data representations across a specific industry sector, group of sectors or across a wider industry base, and develop "standards" for both the application access, data storage layers and interfaces for distributed web-services. Figure 5 shows the basic principle of the data warehouse domain in conceptual terms for ISPDM.



One way to develop the access layer is to use the emerging concept of *business objects* that will be discussed later in section 2.2, System Architecture. The Business Object Application Access layer is designed to simplify access to the database and allow the developers of third party applications to write calls to the database without needing to know its detailed underlying data structures.

**The guiding principals for data management and exchange in the 21<sup>st</sup> century is the adoption of the data warehouse domain and standards.**

## 2. IMPLEMENTATION CHOICES FOR A 21<sup>ST</sup> CENTURY DATA MANAGEMENT SYSTEM

The factors that must be considered in any implementation are:-

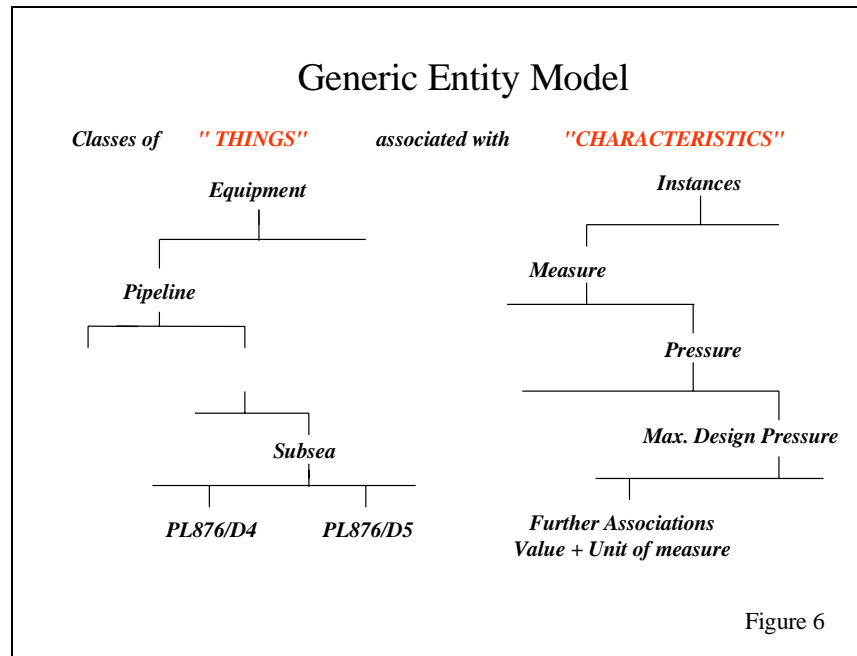
- The Data Model.
- The System Architecture
- The Technologies

### 2.1 The Data Model

The choice of data model is crucial to the success and take-up of any implementation. As noted in 1.1, most previous systems had bespoke underlying data models designed by, or for, specific organisations and applications with data highly coupled. ISPDM is designed to overcome the problems associated with such data models by adopting a “standards” based approach. However, it is important to note that a true standard must be widely agreed by an open process and achieve wide spread use and implementation. Therefore, standards' technology is not static but is a new emerging group of technologies.

As ISPDM is aimed at the pipeline industry, a standard engineering domain model was the logical choice. ISPDM is based on the EPISTLE Core Model, Version 4 (ECM V4) which is

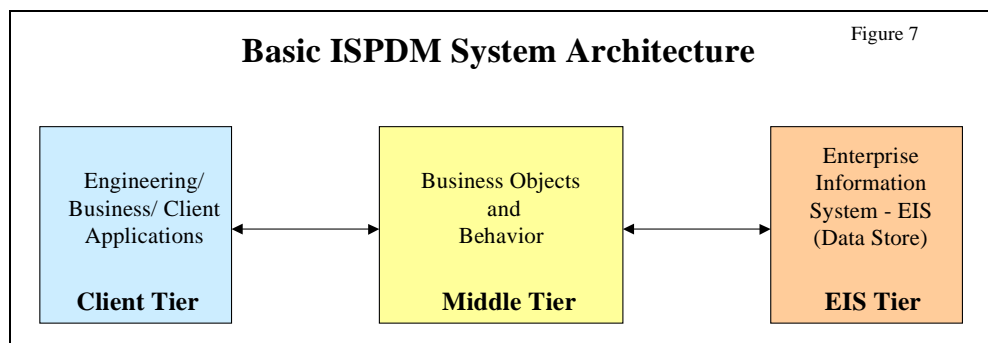
being standardised as ISO 15926 part 2<sup>1</sup>. The basic concept of this model is shown in figure 6.



A challenge for ISPDM was to extend the use of ECM V4 beyond engineering data to include other information required for pipeline asset, risk and integrity management. This required business, environmental, geological, geotechnical, hydrographic and topographic data to be modelled and implemented. If this were achieved successfully it would demonstrate that the EPISTLE model was truly a generic data model that could be applied in major industries and sectors. At the time of writing no major difficulties have been identified in creating an ISO compliant database to hold this wide variety of data.

## 2.2 System Architect

A three tier system architecture was chosen at the outset of the project as shown in figure 7.



<sup>1</sup> ISO 15926 – Integration of life-cycle data for process plants including oil and gas production facilities.

The is a very common design pattern that allows

- Centralisation of Business Logic.
- “Thin client” support.
- Flexibility in the choice of EIS (storage layer).

### 2.2.1 Client Tier

This is the application level. There can be many applications as outlined in figure 5. Hydrographic applications could include paper and electronic chart compilation, updating the generation of lists of lights, wrecks, tidal data, etc., terrain modelling, 3D visualisation and GIS spatial browsing.

The system client tier also supports applications that perform data management functions such as data loading, unloading, updating, deleting, integrity checking, etc. These do not contain any technical or analytical functions outlined above.

### 2.2.2 Middle Tier

It is at this level that the data are assigned to, and merged in, the correct business object and physical storage location in the EIS. It is also the level at which dynamic data sets are derived by data query are built and retrieval from the EIS and forwarded to applications. Business Objects are central to this process. *What are Business Objects?* Business objects are simply:-

- Represent things, processes or events that are meaningful to the conduct of business.
- Provide a way of managing complexity, giving an application perspective.
- Package essential characteristics more completely.
- They can be "industry standard", "company specific", "process specific", etc.

In essence business objects model the behaviour of the data and:-

- Enable users to easily identify the data they require.
- Enable rapidly building of interfaces to link existing applications to the database.

### 2.2.3 EIS Tier: The Data Store

This is the underlying database for holding persistent data.. It is also the area that has proved most problematical, to date, for the project. The problems have been related to the complexity of the “implicit”<sup>2</sup> and highly normalised<sup>3</sup> EECM V4 model with the available database technologies.

---

<sup>2</sup> In an implicit data model much of the data model is held as data in the model schema in comparison to explicit data model where the model is the model schema.

<sup>3</sup> In a highly normalised model data about an entity are mostly represented by relations to classes and entities, and not attached directly to the entity as attributes.



## 2.3 The Technologies

A technology trial always carries a high commercial risk, which is why the EC fund projects such as ISPDM. Therefore by implication ISPDM must break new ground in the application of new and emerging technologies.

### 2.3.1 System Architecture Technologies

The JavaSoft J2EE specification<sup>4</sup> was chosen as the core technology to use for the 3-tier architecture implementation. J2EE is an open standard based on the Java programming language and therefore is portable to a large number of computing platforms. Many commercial implementations of the application-server facilities required in a middle tier are available.

A high level vision of the flexibility this architecture is shown in figure 8. An early “proof of concept” system was built around a small set of sample data that made up Use Case “A”. This subset is illustrated in figure 9. Figure 9 also illustrates that key components from each of the tiers in the multi-tier architecture that have been implemented in order to validate the architecture.

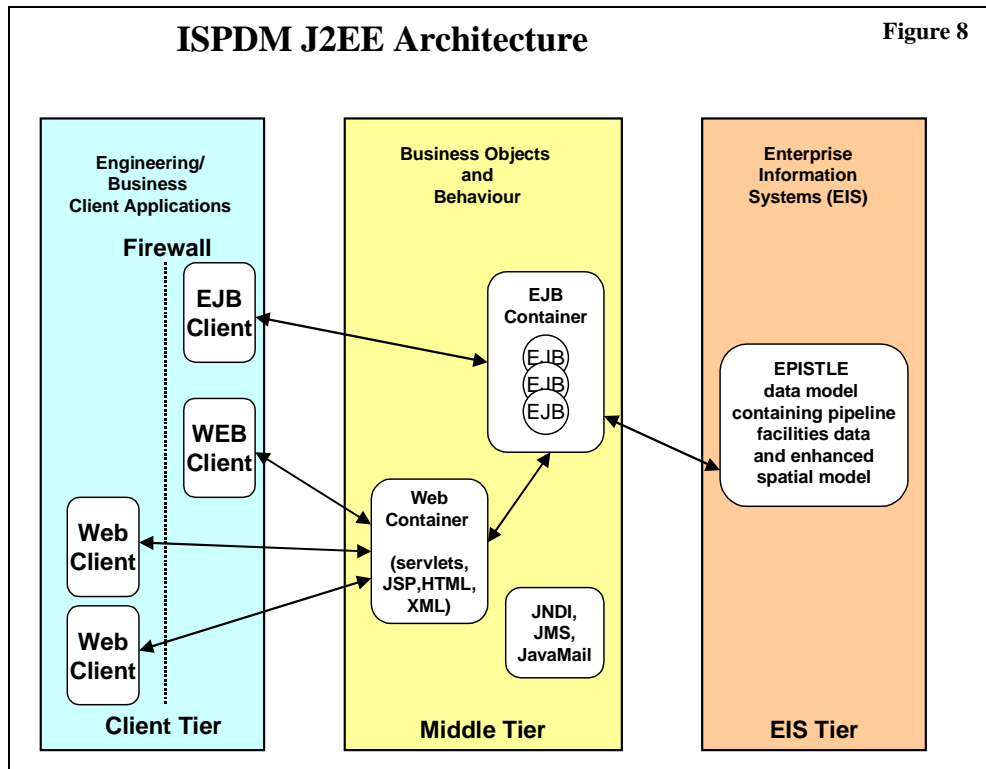
Microsoft’s Internet Explorer V 5.5 has been used for the client tier. This is supplemented with the SVG (Scalable Vector Graphics) plug-in from Adobe. The middle tier uses Borland AppServer V4.1 to provide a web container. The web container processes the Java Server Pages requested by the Client Tier and uses a Server Bean to contain the business logic. The EIS tier uses the PrismTech Synergy cartridge (at present in a fairly minimal way) together with the file system. The Synergy<sup>5</sup> cartridge is an implementation of ECM3<sup>6</sup>. The current output, generated by the JSP and bean components in the middle tier can be viewed in a web browser.

---

<sup>4</sup> The Java 2 Platform Enterprise Edition (J2EE) architecture defined by Sun Microsystems Inc. for implementing enterprise applications.

<sup>5</sup> Synergy Cartridge; Developed for the Synergy project by PrismTech as an implementation of EPISTLE in Oracle and originally intended to be bundled with Oracle. This bundling has not happened.

<sup>6</sup> EPISTLE Core Model version 3



After proving that these technologies worked, the issue of potential users who preferred a Windows environment was addressed. The SOAP<sup>7</sup>, Simple Object Access Protocol, is just emerging. As with J2EE the SOAP environment is platform independent and relies on XML to define the format of the information and then adds the necessary HTTP headers to send it. A test was performed to assess the use of this technology in relation to ISPDM. In a little over two weeks a SOAP architecture was put together and data loaded, retrieved and displayed in a beta release of Intergraph's GeoMedia web map product.

This architecture is shown in figure 10. As it shows issues associated with GIS that will be discussed shortly. At time of writing the project is reviewing its capability, within budget constraints, to develop a prototype with both J2EE and SOAP capability.

<sup>7</sup> SOAP= Developed by Microsoft, DevelopMentor & Userland Software.

Figure 9

### ISPDM First Iteration J2EE Implementation - USE CASE A

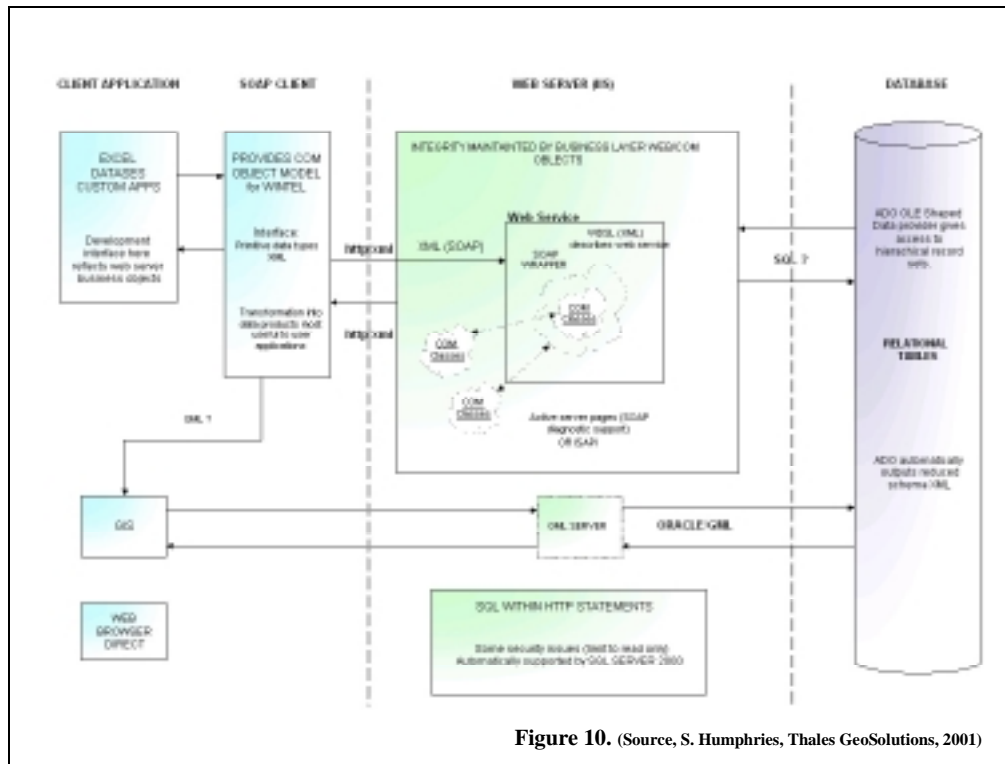
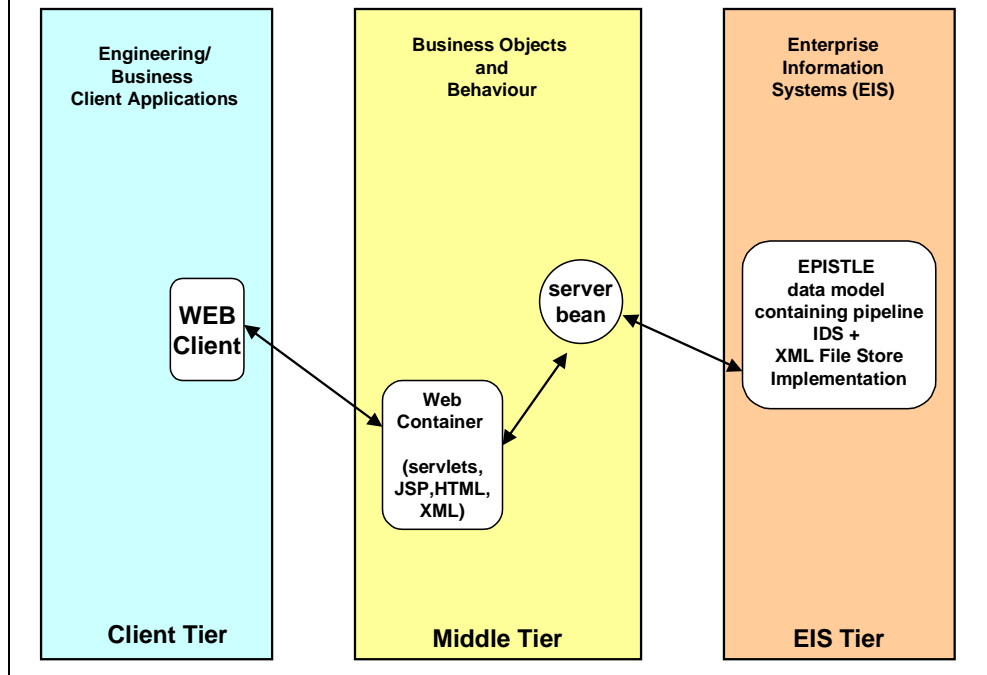
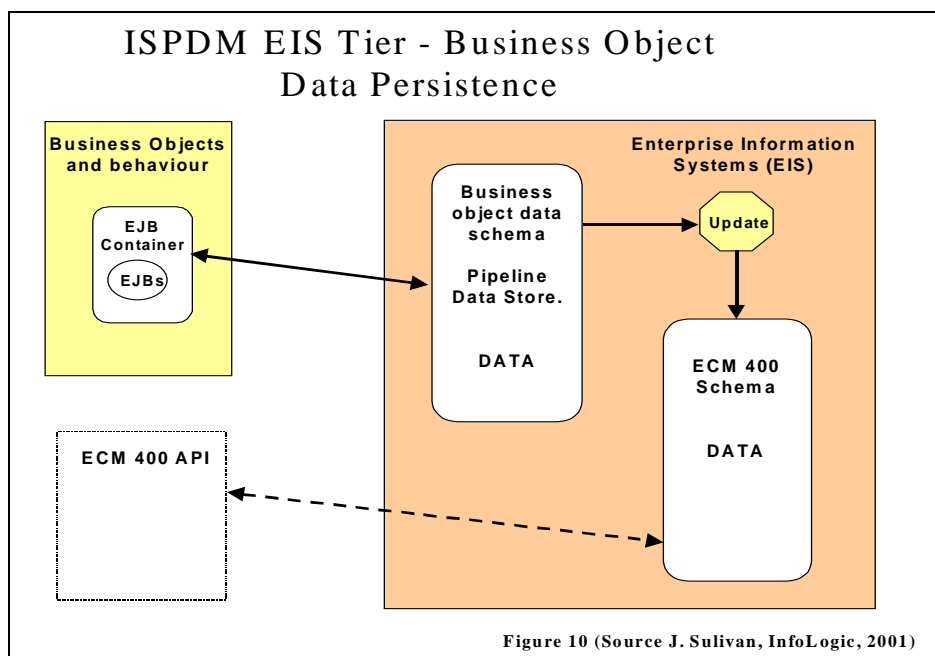


Figure 10. (Source, S. Humphries, Thales GeoSolutions, 2001)

### 2.3.2 The EIS Storage Layer

As noted earlier this proved to be the most problematical. A drawback with ECM V4 is that it does not deliver the retrieval performance that most users of proprietary databases have come to expect. This is due to it being an implicit normalised model where entity attribute data are stored a separate entities linked by relations. The Synergy Cartridge technology was expected to overcome this problem. However the Synergy cartridge was designed for ECM Version 3. However, ECM Version 4<sup>8</sup> is considerably different and re-engineering the cartridge to be compatible was a cost the project could not bear.

Instead of re-developing the Synergy Cartridge, ISPDM chose to develop a new EIS business object tier based on figure 11. The EIS data store will contain two levels. The first is the data held in the form of business objects and the second the ECM4 compliant data store. For normal use business objects will be accessed and the data passed back to the applications through the middle tier. The ECM4 data store will act more as an archive. A useful by-product of this is that the data store could be accessed directly by an ECM API for retrieval or exchange of data in ISO 15926 part 21 files. A number of major E&P companies now make this a mandatory requirement for major database products.



### 2.3.3 Communications and XML

XML<sup>9</sup>, Extensible Markup Language, is the emerging universal format for structured documents and data on the web. As ISPDM was intended to be web enabled from the outset, XML was the obvious choice for passing information between the various tiers in the system.

<sup>8</sup> EPISTLE Core Model version 4.

<sup>9</sup> XML is being developed under the umbrella of the World Wide Web Consortium (W3C).

There are a number of XML software packages on the market. The project is using XML-Spy 4.0 for the generation of the XML Schema and the XML document.

GML<sup>10</sup>, Geography Markup Language, was chosen as the mechanism to link GIS applications, and SVG to link web browsers, to the system. GML is an XML encoding for the transport and storage of geographic information. ISPDM will not use GML in this way, as it cannot be made to conform to the ISO 15926 standard. Instead GML will be generated from ISPDM's XML documents. The links to GIS and SVG applications have been tested and proven.

### **3. THE PROCESS OF CREATING A 21<sup>ST</sup> CENTURY DATA MANAGEMENT SYSTEM**

The process of implementing ISPDM was as follows:-

- Define the data types required to be stored.
- Identify and define the business objects.
- Model the business objects in UML to be ECM4 compliant.
- Create the XML document definitions.
- Build an Oracle ECM4 pipeline data store.
- Create data loading and retrieval mechanism.
- Create an underlying full ECM4 data store.

#### **3.1 Define the ISPDM Data Types**

This work was an extension of part of the project feasibility study conducted in 1999. The emphasis was to define all the data and information that an engineer would require to conduct a comprehensive integrity management and risk assessment analysis of a pipeline.

Typical information includes engineering design data, operational data, environmental data, inspection and maintenance data, hydrographic data, topographic data and general business data. Data of particular interest to this audience includes:-

- Environmental data, including wave, wind, current and temperature data.
- Geotechnical data, including seabed conditions and soils properties.
- Hydrographic data, including pipeline risk information such as restricted areas, shipping lanes, fishing activity.

To date some 1715 types of information have been identified. However, many the data types are used more than once and 3363 data fields have been identified. For example only one "easting" is included in the 1715 while there are 220 instances of "easting" in the 3363 data fields. A small sample of this is included in box 3.1 These data are detailed in a project document<sup>11</sup>.

---

<sup>10</sup> GML is being developed by the Open GIS Consortium

<sup>11</sup> ISPDM Data Types Document V5\_Rev 05\_7.01.02

Data Item	Occurs In Data Types Document V5 - Data Entry
Easting	End of reinstated area (Reinstatement – Basic installation & construction data)
Easting	End of section (Noise abatement - Basic installation & construction data)
Easting	End of trench (Trenching – Basic installation & construction data)
Easting	From (Operations - Basic installation & construction data)
Easting	Intermediate points (Trenching - Basic installation & construction data)
Easting	Location (Field bending – Bends - Basic installation & construction data)
Easting	Location (Pups - Basic installation & construction data)
Easting	Location (Tie-ins - Basic installation & construction data)
Easting	Start of reinstated area (Reinstatement – Basic installation & construction data)
Easting	Start of section (Noise abatement – Basic installation & construction data)
Easting	Start of trench (Trenching – Basic installation & construction data)
Easting	KPO
Easting	Location(s) (Additional stresses - Wall thickness design - design data)
Easting	Location(s) (Collapse - Wall thickness design – design data)
Easting	Location(s) (Critical spans – design data)
Easting	Location(s) (current data – design current data)

Box 3.1

### 3.2 Identify and Define the Business Objects.

From the data types document version 5 some 210 business objects have been identified. These can be grouped into about 25 business object types. A sample of the business objects and their type grouping for objects with a hydrographic content are shown in box 3.2.

Business Objects including Hierarchy		Business Object Definition
Individual	EnvironmentalElement	
	Atmosphere	Wind
	MarineGrowthArea	
	SeabedArea	SoilSeabed
	WaterArea	RockSeabed
		WaterCurrent
		WaterTide
		WaterWave
	DesignatedBoundaries	
	LicenceBlockBoundary	
	InternationalBoundary	
	OtherBoundary	
	DesignatedArea	
	LicenceBlock	
	LandParcel	
RestrictedArea		
FishingArea		
		Gaseous envelope surrounding the earth.
		Air in more or less rapidly moving motion.
		An area of seabed with marine animal or plant life attached to it.
		An area of seabed with common soil properties.
		An area of seabed mainly composed of rock.
		A spatial temporal part of seas or lakes
		Water moving in a given direction
		Periodic rise and fall of sea level due to attraction of moon and sun.
		Ridge of water between two depressions. (Oxford Dictionary)
		Formally defined, and possibly legally defined, boundaries.
		Boundary defining an E&P licence area.
		Boundary between two sovereign countries that may or may not be ratified.
		Boundary defining a generally recognised area that may or may not be legally defined. E.g. Territorial Sea, EEZ.
		An area with some common usage characteristics that may or may not be legally defined.
		Area formally designated and defined for the licensing and control of E&P activities.
		An area of land designated by ownership.
		Area with certain characteristics that result in the formally recognised restriction of activities not associated with the characteristics. E.g. Military Practice Area, Nature Reserve, Spawning Grounds, etc.
		An area of sea that is an established or recognised fishing ground.

Box 3.2

An integral part of the standard is the ERDL, EPISTLE Reference Data Library, which defines and describes the terms used in the data model. This is published as ISO 15926, part 4. The final definition included in the ERDL becomes the accepted ISO definition of that data item. Thus any hydrographic features included in ISPDM and defined in the ERDL will have de-facto ISO definitions.

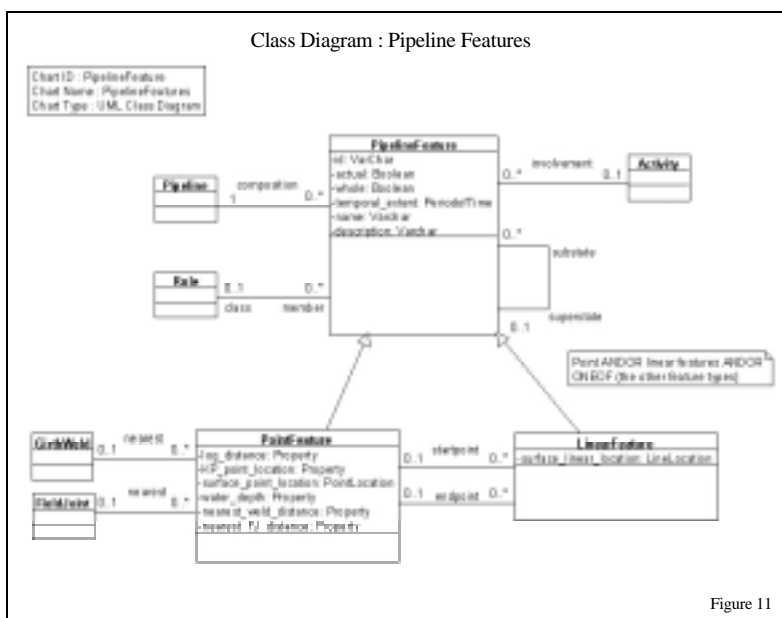
It is interesting to note that the project has had to create the business objects WaterTide, WaterCurrent and WaterWave instead of using the traditional hydrographic terms Tide, Current and Wave. This is because that can be no ambiguity in a data standard. There are other types of tides, currents and waves such as gravimetric tides and electrical currents and radio waves.

### 3.3 UML Modelling of Business Objects

The data modelling of the business objects is a very crucial element of the process. If this is not correct then the data will be misrepresented in domain terms and stored in an inappropriate and non-compliant manner. This required detailed understanding of the model and, or access to, domain knowledge. The UML modelling language was chosen for the business object model.

The UML classes defined in this process are always specialisations of the more general concepts defined by the ECM V4 model. The chosen specialisations are obviously those that are useful in the pipeline domain. This rule enables the business object tier to be mapped to the ECM V4 compliant tier.

The project is fortunate is having the services of a data modeller with intimate knowledge of ECM V4. The individual has played a central role in the model's design and development. The consortium members provide extensive domain knowledge and by close liaison between all parties valid UML diagrams are produced. The importance of this interaction between data modeller and domain experts cannot be overstated.



At the time of writing UML modelling has not been done for hydrographic data. It has been done for a sample of engineering data. The following figures are examples of the UML class diagrams. Figure 11 is the generalised diagram for pipeline feature. Figures 12 and 13 are the UML diagrams for two types of features; a span, an external feature, and a metal loss feature, an internal corrosion feature.

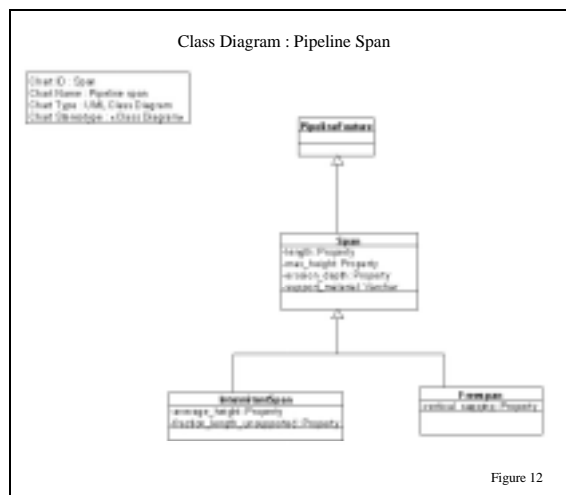


Figure 12

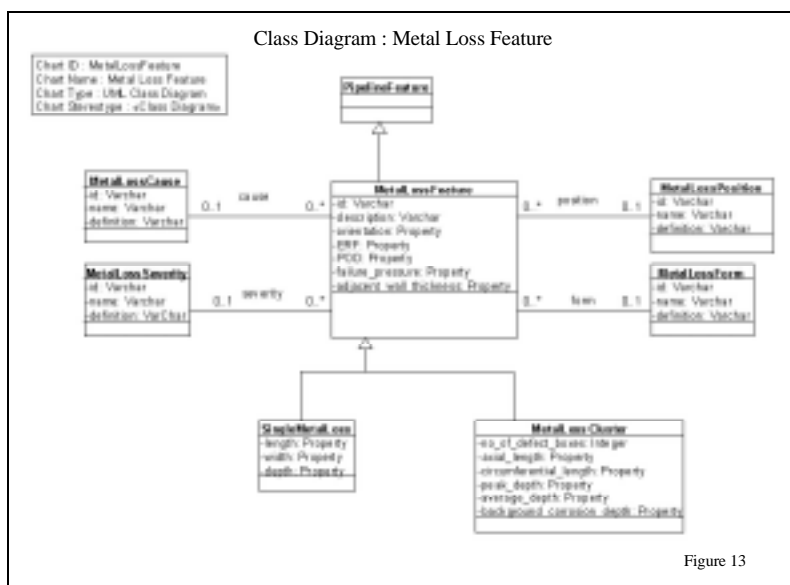


Figure 13

### 3.4 The XML Document

The first “proof of concept” project iteration involved the generation of a DTD, Document Type Definition, and an XML document using XML-SPY. However, subsequent deliberations concluded that this was inappropriate. The DTD was replaced by an XML Schema that better reflected the UML diagrams and hence the ECM4, ISO, data model. The XML documents were then constructed from the XML Schema.

The process of UML modelling and the generation of the XML Schema means that the resulting XML documents conform to the tenets of ECM4. Therefore, the XML document,



used for both transport and storage of data, are ISO 15926 compliant. As a result, once defined, the XML documents can only be extended and amended by following the rigorous process outlined above. This is a somewhat unique approach to the construction and use of XML.

For those not familiar with XML, figures 14 to 17 are taken from the XML editor, XML-Spy and show the construction of the document. Figure 14 shows the document at its highest level. Figures 15 to 17 show it progressively expanded until, in figure 17, some data can be seen.



Figure 14

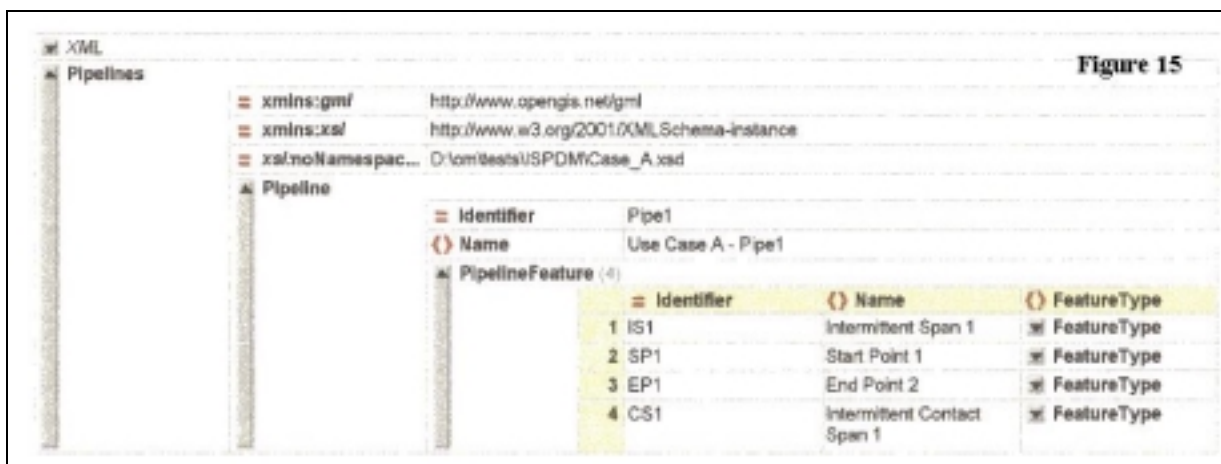


Figure 15

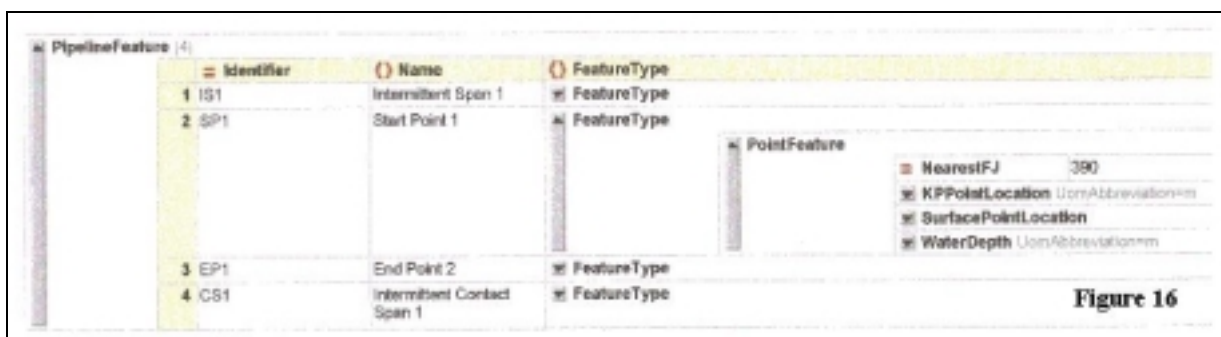


Figure 16

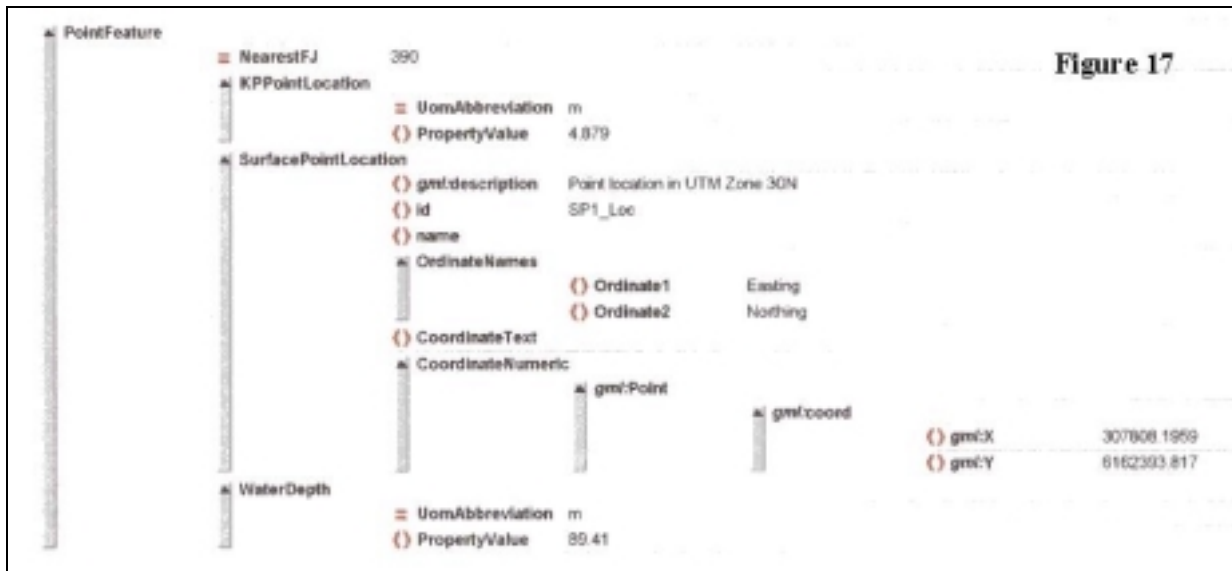


Figure 17

Figure 18, is a section of the actual XML document that will be transported over the Web which contains the data. The section of the document within the box contains the information shown in figure 17.

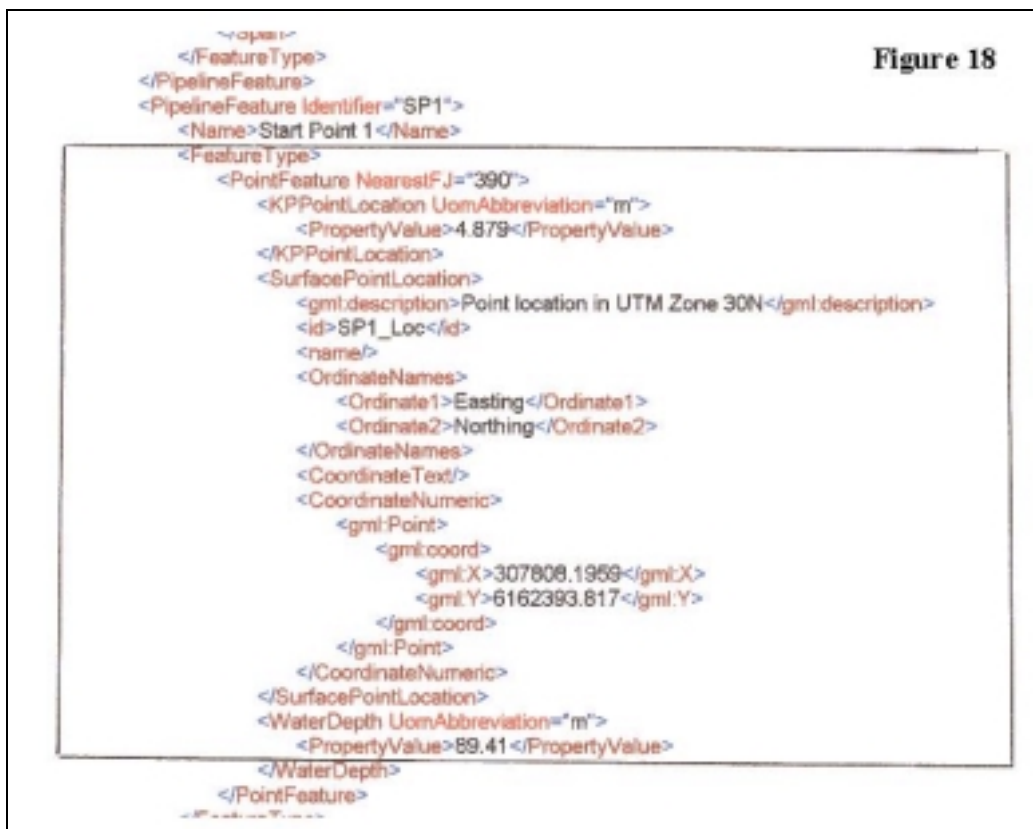


Figure 18

### **3.5 The Underlying Data Store**

Design of the underlying data store was a significant early problem for the project. Apart from the Synergy cartridge failing there was the issue of how to handle spatial data. This is crucial to allow GIS applications to operate against the database system.

Initially it was intended to use Oracle 8i and duplicate the spatial data within the data store. Coordinates would be held in their original form in the ECM4 data store and duplicated in the Oracle spatial Cartridge, OSC. GIS applications would directly link to the OSC to perform spatial queries. However, with the advent of Oracle 9 this duplication appears unlikely. The separate functionality of the OSC is now integrated in Oracle 9. At the time of writing the project is evaluating the use of Oracle 9 and its potential to solve the spatial problem.

The tables for the “Pipeline” data store, shown in figure 11 are developed from the UML schema to hold business object data. Another set of tables holds the data according to the ECM V4 model. In this way the data store is fully ISO 15926 compliant.

### **3.6 Data Loading and Retrieval**

The use of XML for data transport had been established at the outset. This requires that applications read to and from XML documents. In terms of the pipeline engineering domain Excel spreadsheets are the basis of many applications. A test was performed early in the project to demonstrate that using VB scripts, XML can be automatically generated for Excel and visa-a-versa. In the near future this facility will be an integral part of the Microsoft office suite of software. What is not known, is how this will be able to be used against the type of structure document being developed by the project.

For loading and retrieving bulk data sets it may be more efficient to bypass the XML and map the bulk file directly to the database tables. However, in an effort to ensure standardisation of data transport the project may opt not to do this but stick with the XML route. In the short to medium term this will create a standardised data exchange mechanism between ISPDM and non-ISO standard system. Links to GIS using GML have already been discussed.

## **4. THE SPATIAL & GIS PARADIGMS**

The EPISTLE model does not contain a geodetic reference framework. It could be modeled and added. However, the project does not intend to duplicate work already in the public domain. Therefore, it is proposed to hold an existing framework, such as the EPSG or NIMA frameworks, within Oracle 9.

The method of handling coordinate information has already been discussed. However ISPDM threw up what has been termed the spatial paradigm, as outlined in figure 19. In terms of the business objects it could be argued that a GIS system would be suitable as a pipeline database. However, when looking at all data types the spatial element is only 9%. Therefore, GIS is clearly not a suitable tool for data management.

A similar assessment was made of the relevant hydrographic business objects incorporated in ISPDM. Although 100% of objects had a spatial component the spatial element only amounted to 17% of the data associated with them (figure 20).

<i>ISPDM Spatial Paradigm</i>		
	<i>Business Objects</i>	<i>Minor Data Types (Characteristics)</i>
<i>Spatial Occurrence Percentage</i>	91%	9%

*Figure 19: The Spatial Paradigm*

<i>Hydrographic Spatial Paradigm</i>		
	<i>Business Objects</i>	<i>Minor Data Types (Characteristics)</i>
<i>Spatial Occurrence Percentage</i>	100%	17%

*Figure 20: The Spatial Paradigm*

## 5. CONCLUSIONS

The ISPDM project clearly demonstrates that the new and emerging standards, data warehousing and web technologies can be successfully implemented in a domain that involves a wide variety of data types, including hydrographic data, and be linked to GIS, mapping and charting, and other graphical applications.

The advantages to hydrographic community of such technology are:-

- *Long term cost savings* - simpler and consistent access to, and exchange of, data will lead to saving both internal and external to any HO.
- *Full life-cycle capability* - the proposed underlying EPISTLE data model incorporates business and full life-cycle, historical, information storage capability.
- *Vendor independence* - a standard frees organisations from the dependency on any one particular software application vendor.
- *Resource maximisation* – vendor independence allows freedom to maximise resources by changing vendors and applications as technology changes.
- *Improved vendor applications* – vendors can concentrate on the enhancement of their applications for data capture, verification, manipulation and ENC output knowing that data will be accessible in a consistent format no matter who the client HO is. This increases the vendors market, encourages competition and drives down the HO's costs.

- *Future proof* – a system based on the EPISTLE standard will largely be technology independent and stable over a long period of time.
- *Model support* – The basic data model will be supported under the ISO umbrella thus freeing the HO from the costly need to continually support a proprietary model.
- *Maritime safety* – The development of web based standard complaint database systems will allow ready exchange of data between HO's and between HO's and ENC vendors thereby negating the need for ENC vendor to maintain independent data sets. The core data will be managed by the competent authorities thus ensuring that only rigorously quality controlled data will be presented in ENC's.

## **BIOGRAPHICAL NOTE**

The author graduated from Glasgow University in 1972 with a degree in Topographic Science. He subsequently obtained his MBA from Glasgow in 1992 and is currently undertaking a part time PhD at Cranfield University. He has spent 30 years in the offshore oil and gas industry and held a number of senior survey related posts in BNOC, Britoil and BP. He has conducted significant research into issues surrounding data management and pipeline data management and application requirements. This research lead to the ISPDM concept. The author currently is project managing this EC funded project. The author has published numerous articles and presentations on issue such as deep water surveys, quality management