

UNICORN SOLUTIONS, INC.
PROPRIETARY & CONFIDENTIAL

CENTRAL INFORMATION MODELS FOR DATA TRANSFORMATION



This document originally appeared in
EAI Journal (www.EAIJournal.com), Fall 2002
by Joshua Fox
www.Unicorn.com

“CENTRAL INFORMATION MODELS FOR DATA TRANSFORMATION”

Regardless of industry specialization or organizational size, enterprises face critical problems in data management and operational efficiency as a direct result of multiple overlapping and distinct schemas for each business domain. Corporate mergers may result in parallel applications that handle similar business domains; in many cases, different departments may define their transmission schemas separately because of a lack of coordination or the use of different applications.

Maintaining the variety of applications and integrating between them is a tremendous challenge. It is difficult to track the schemas and the business entities that each one represents. At present, to integrate the applications that use these schemas, transformation logic must be custom-developed at great cost, reflecting semantic and syntactic differences between schemas. Any effort to manage these metadata and data-transformation assets manually is doomed to failure, since changing requirements quickly invalidate analyses and inter-schema transformations.

THE PROBLEM

Enterprise data resources can take many different forms: relational database tables, XML documents, EDI messages, COBOL records and others. Independent Software Vendor (ISV) applications such as enterprise resource planning (ERP), customer relationship management (CRM) software, or home-brew software, each defines its own schemas for input and output. In many cases, the different schemas hold similar information, although the structure of the information may differ.

EAI systems must overcome the problem of heterogeneous schemas. To do this, EAI implementers must fully understand the business semantics of each schema, since the schemas often represent the decisions of a single application developer working from specifications of a single step in a business process. Over time, new schemas multiply as business requirements change, and the investment in analyzing business requirements, as captured in older schemas, is lost.

Today, overcoming this heterogeneity requires significant investment of time, with little return. Developers must create transformation code manually, whether by coding in textual languages such as XSLT, or by using graphical transformation design tools. While this is tenable with a few schemas, a larger number of schemas results in a skyrocketing number of such transformations. The number of transformations goes up as the square of the number of schemas.

Even if applications are manually integrated at great expense, the solution is not maintainable. When schemas change, as they often do in a quickly changing business environment, existing code for transforming data from one schema to several others becomes unusable. In the worst-case scenario, when a single application's output or input schema changes, all code for integrating that application must be rewritten from scratch, since XSLT code and SQL queries do not lend themselves to reuse. And since the business semantics of the original schema are frequently not captured formally, re-developing transformation logic means re-analyzing the semantics of the schemas. Finding the appropriate domain expert and dedicating the time required for this additional development cycle are just two of the obvious drawbacks.

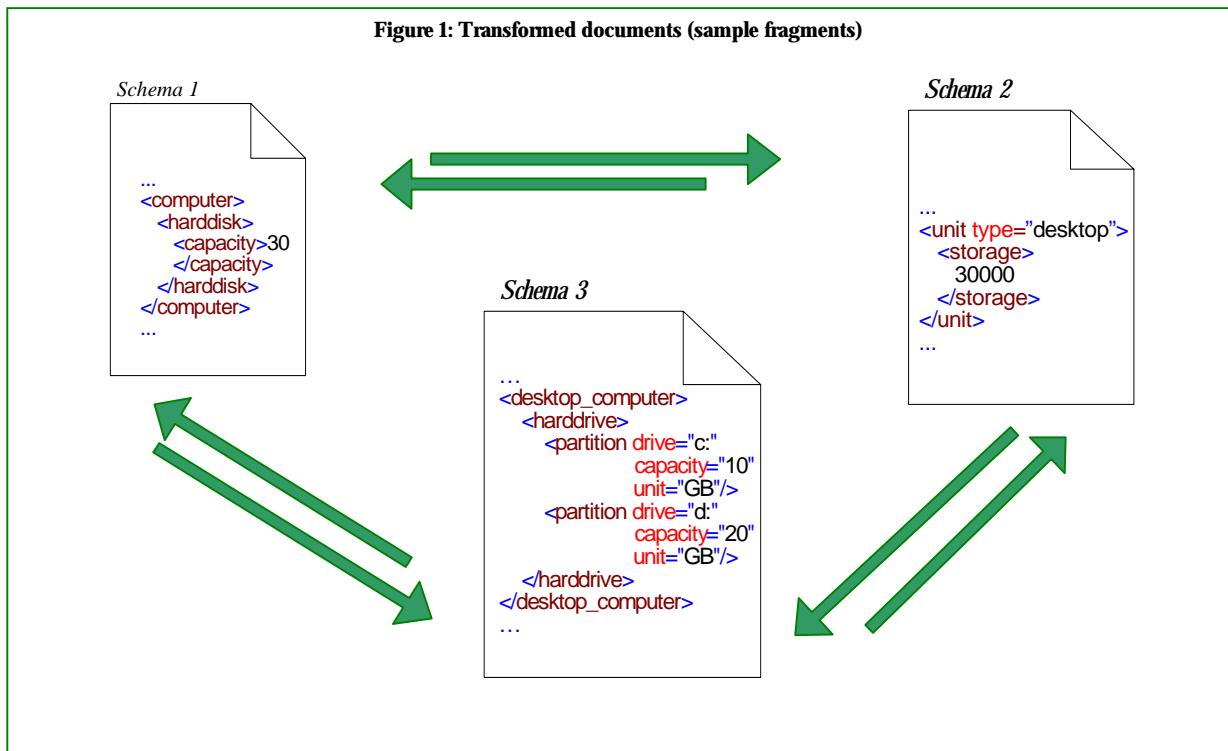
To illustrate with a simple example from the area of high-technology manufacturing (see document fragments in Figure 1), desktop computers with their associated storage capacities may be represented with a variety of XML structures. Each schema conveys fundamentally the same information (a desktop computer with 30 gigabytes of hard-disk capacity) but the tag names and structure are completely different. In our scenario, each of these applications must send messages to the others on the bus. To convert from one to the other, elements must be taken out of the

source document and arranged in a different way in the target, with some application of business rules along the way: for example, converting the megabytes recorded in Schema 2 to the gigabytes of Schemas 1. If the three applications here are to be fully integrated, six XSLT transformations will have to be written.

Although the example here is built with XML, the same principles would apply to other structures such as EDI messages, or even to disparate database tables, transformed by means of SQL.

TODAY'S EAI SYSTEMS

Today's hub-and-spoke messaging bus architectures have revolutionized application integration, addressing some of the problems posed by point-to-point heterogeneous enterprise



information systems (See Figure 2). However, they still require point-to-point effort for schema integration.

Figure 2: Hub-and-spoke on architectural layers	
Integration Layer	Hub
Transport	EAI bus
Format	XML
Schema	Central Information Model

The *transport* hub-and-spoke architecture, as represented by asynchronous messaging buses, has simplified integration efforts tremendously: No longer do integrators have to develop a messaging link for each pair of source and target applications; the source application can now send a message on the bus, while the target application simply listens to an agreed-on queue or "topic."

The *format* hub-and-spoke architecture deals with the wide variety of data formats found in the enterprise, such as comma-separated textual fields, application-specific binary protocols, and relational databases. A standards-based approach uses XML as a hub format. EAI technologies can take data from most common formats and convert it to the hub format. Nonetheless, this solution to the problem of disparate *formats* does not solve the problem of disparate *schemas*. Existing tools can convert all messages to XML, but the schemas may differ. Similar information may be carried in different elements or attributes, as in our example (Figure 1).

EAI vendors include message broker components (also known as “integrators”) that can transform data from one schema to another, once transformation logic is deployed. These message brokers sit as a peer application on the bus. Each source application sends messages to a messaging-bus address for the message broker. The message broker receives the message, and, depending on workflow rules, transforms the message into a schema suitable for a target application, and then re-sends the output to the target. These message-broker transformation engines are an essential part of modern EAI systems.

However, the message brokers still require the transformation logic to be developed manually. Thus, a *schema* hub-and-spoke architecture is still needed to avoid the design-time effort of application-to-application transformation.

EXISTING TECHNIQUES FOR TRANSFORMATION DEVELOPMENT

Enterprises today have various approaches to the problem of multiple incompatible schemas. The typical solution is to hand-code distributed services that transform one schema to another, using a general-purpose language such as Java, or a transformation-specific language such as XSLT or IBM’s ESQL.

Another approach is to develop transformation logic using graphical development tools, often bundled with EAI message broker products. Although these tools simplify the development of any given transformation, the user must still create transformation logic manually for a given source schema and target schema. This includes analyzing the business-requirements of each schema, designing the transformations, and building them. The effort required still goes up rapidly as the number of schemas increases, and it is still impossible to maintain the transformation logic in the face of application requirements. This is parallel to the well-known problem with application-to-application transport connections that was solved by messaging busses: an integration effort that goes up rapidly as the enterprise’s business needs inevitably grow.

THE UNIFYING INFORMATION MODEL

The semantic hub is at the heart of a new enterprise data integration architecture in which a central information model dynamically generates data transformations between and among schemas, obviating the need to develop schema-to-schema transformation code for any pair of schemas.

This new generation of software creates a new hub-and-spoke system for *schemas*, naturally extending the benefits of the hub-and-spoke system for *transport* and for *format* found in today’s EAI systems. This semantic hub automatically generates all the needed transformation code.

The semantic hub is centered on a rich, central, active information model. This information model is based on ontology, a modeling technique that has been under development in academia for decades. In recent years, Tim Berners-Lee, inventor of the World Wide Web, has put forward his vision of the Semantic Web, the next stage of information sharing, this time between

applications. Ontology is at the core of the Semantic Web. The World Wide Web Consortium, which has created standards for the Web and for XML, is now finalizing a standard approach to ontology as part of its Web Ontology Working Group, guaranteeing widespread standardization for semantic systems. The semantic hub, as used in EAI systems, is an application of the Semantic Web concepts to the particular needs of the enterprise, which include robustness, maintainability, and scalability.

Older modeling techniques have some superficial resemblance to the ontological method. Ontology resembles the Entity Relationship (ER) model in the linking between classes/entities by means of properties/relationships; it resembles Object Oriented (OO) design in the use of inheritance as a powerful class-building principle. However, ER is intended primarily for modeling relational databases, while OO is aimed at software design. Ontology, on the other hand, models the semantics of the enterprise, rationalizing schemas written in a variety of languages, including XML DTDs or schemas, COBOL copybooks, RDB data dictionaries, and others.

The Information Model: Rich

This information model is rich in that it describes real-world business entities, using techniques such as inheritance and inter-class links to indicate specializations among entity types. It includes business rules that show the relationships between business entities. The model presents a coherent view of the complex web of meaning relating business entities in the enterprise. It uses vocabulary from the business domain rather than the cryptic identifiers found in code.

In addition to code generation, the rich information model maximizes the business value of schemas by serving as a central semantic repository for enterprise schemas. Centralizing and clarifying the semantics of a schema in a unified, coherent way, the information model allows business managers to catalog and assess their metadata assets. Legacy schemas are often difficult to understand, maintain, and modify. When schemas are rationalized to the central information model, which uses intuitive business terminology rather than the cryptic identifiers of schema languages, the meaning of the schemas and their sub-elements is clearly indicated.

The active information model maintains full synchronization with enterprise schemas. Mapping each new schema into the model immediately relates it semantically to every other schema. When changes are required by new business requirements, the user can introduce mappings of the new schemas, or the information model can actually generate new schemas.

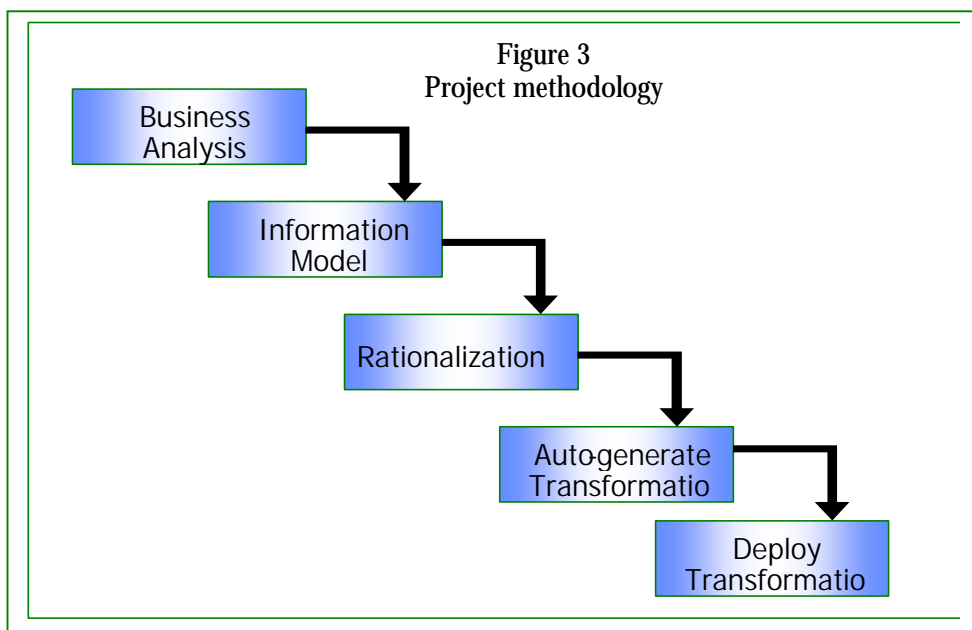
The Information Model: Central

The information model is central in that it represents a neutral semantic view of the enterprise, not related to any one schema. This model is linked to metadata such as XML schemas, Java APIs, COBOL copybooks, and RDB data dictionaries. This is in contrast to other types of models, such as Entity-Relational diagrams or UML, which are related only to the databases or software that they served to design; E-R and UML models are often created during development and then filed away. A central information model, on the other hand, gives coherence to the enterprise view, providing a live picture of the current state of the wide variety of enterprise data resources. Each of the resources is rationalized through mapping to the central model.

The Information Model: Active

The information model is active in that it can generate transformation code to transform from any one schema to any other, logically unifying all schemas. This is essential in application integration, since when two applications must communicate, the output of one differs in its structure from the expected input of the other. The semantic software can generate the transformation logic, seamlessly integrating the two applications without further human intervention.

Creation of this rich, central, active model gives IT managers new power. It converts *data*, comprehensible only to the applications that read and write it, into *information*, which includes a semantic description that allows all enterprise applications to use it.



DATA INTEGRATION PROJECT METHODOLOGY WITH THE ACTIVE INFORMATION MODEL

The development process for creating transformations through an active information model differs from the process for developing them manually: by focusing on schema rationalization to a central information model, the investment in developer time is not lost with the creation of each new transformation. Rather, a single per-schema effort translates into a rich central model that rationalizes and clarifies enterprise data. Transformations can then be produced with no additional developer effort.

This structured process, centered on the information model, minimizes human errors, allows applications to be integrated with greater efficiency, and reduces lag time in responding to changing business needs. The stages of the process include business analysis, creation of an information model, rationalization of the schemas through mapping to the model, automated generation of transformations, and finally, deployment.

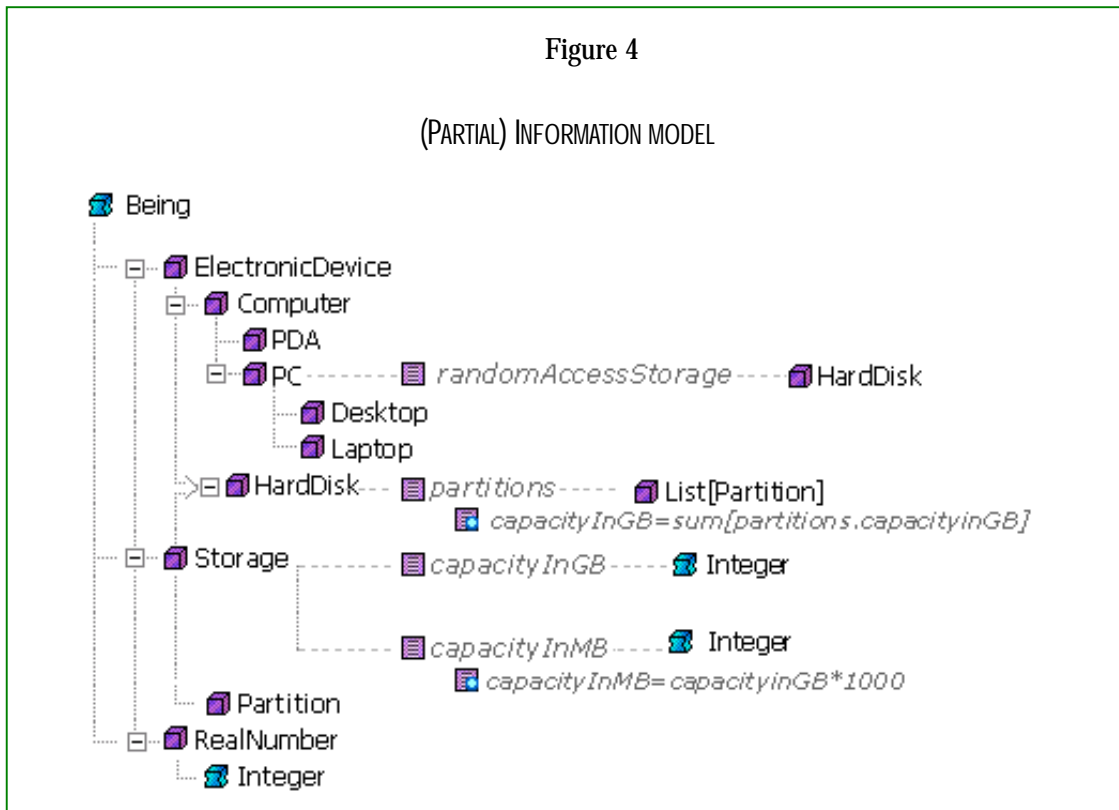
(1) Business analysis

Developing transformations with a semantic hub begins with an examination of data schemas in the enterprise by business domain analysts. The analysts determine the semantic value of the schema elements—the real-world entities that they represent—as well as the business rules that relate the schemas. In our manufacturing example, business analysts determine the computer types being sold, the properties of the computer most relevant to the manufacturing and sales process, the units of measurement for these properties, and the business rules relating the properties. For example, computers are one of the products being manufactured, drives can be divided into partitions, and so on.

(2) Central information model

The second stage is the development of a central information model which encodes the semantic information in the schemas. A small part of a central information model for our computer-manufacturing example is shown in Figure 4. In this example, the class called **ElectronicDevice** is among those that inherit directly from the universal base class **Being**. **Computer** and **HardDisk** are specializations (subclasses) of **ElectronicDevice**; **HardDisk** is also a subclass of **Storage**.

Examining the class **Computer**, we note that each computer has a *randomAccessStorage* device of type **HardDisk**. **HardDisk** has the property *partitions*, showing that the disk has multiple partitions. The information model also has business rules, for example, a rule indicating that the capacity in megabytes equals the capacity in gigabytes multiplied by 1000.



(3) Rationalizing the schemas

The third stage of the development process is to rationalize the schemas by mapping them to the ontology. In this stage, each complex (structured) element of the schemas is mapped to a class, while each simple (atomic) element is mapped to a property of a class.

In our example, the structured elements representing a desktop computer in each of the three schemas are mapped to class **DesktopComputer**; the simple elements indicating capacity in megabytes are mapped to the property *capacityInMB*.

At this stage, the advantage of the semantic hub becomes clear. The schema-to-schema approach requires the development of transformation code for each pair of schemas, while the semantic hub approach requires each schema to be mapped only once to the central information model.

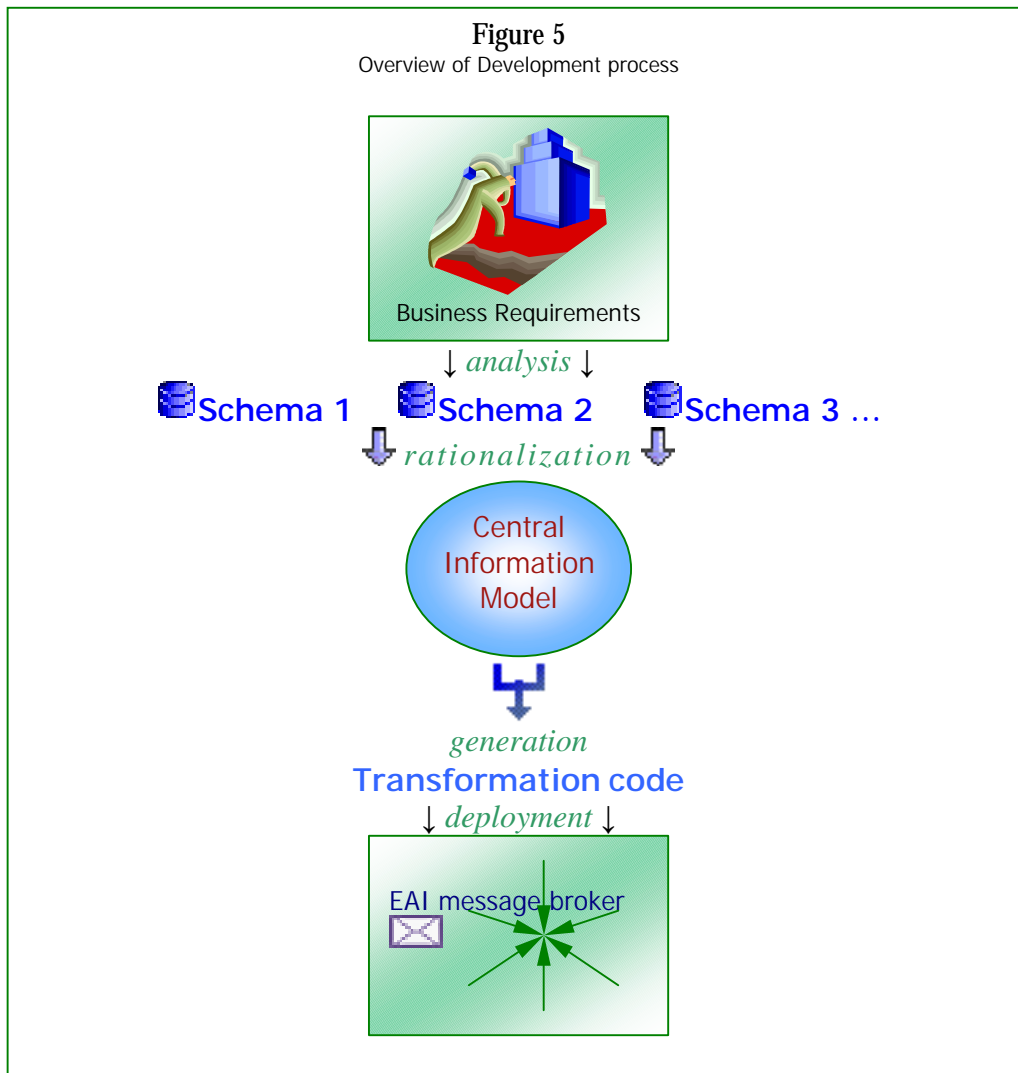
(4) Auto-generating transformation code

The fourth stage is the automated generation of transformation code. The semantic hub software searches for shared semantics of the source and target schemas, where both source and target are mapped to the same concepts in the model. More advanced semantic hubs can transform schemas even if they are not mapped to the same concepts in the model, by applying business rules encoded in the semantic hub.

In our example, the semantic hub encodes the fact that different elements are mapped to the same class: **computer**, **unit**, and **desktop_computer** elements are all mapped to the **Desktop** class. When converting XML from one schema to another, the respective subdocuments are transferred to the new tag names. Likewise, the multiple elements called **partition** in Schema 3 are transformed into a single **harddisk** element in Schema 1. The semantic hub further generates the code needed to convert megabytes into gigabytes under the encoded business rule.

The generated code can be in any of a variety of languages: for XML, as used in EAI messaging systems, XSLT is used. When the goal is to directly transform between relational databases, the code is SQL that SELECTs data from the source RDB, then INSERTs the transformed data into the target.

Auto-generated code gives the additional advantage of maintainability: As business needs change, a simple update to the active information model allows for new code generation.



(5) Deployment

The final stage is deployment. Just like manually developed schema-to-schema transformations, the auto-generated transformations can be deployed into the message broker component of an EAI system, then executed in the message broker's transformation engine.

CONCLUSION

Enterprises are faced with too much data, too many data schemas, and not enough information. Semantics—meaning—is what turns data into information, but often the semantics are expressed only implicitly in application code, so that application data interchange requires human input.

Enterprises are now undergoing a slow revolution from managing data to managing information. For data to go beyond the silo walls of a single application to full cross-application sharing in a smooth and automatic manner, semantics must be formally expressed and then linked to the data. A rich information model makes data comprehensible to external applications designed to read it. EAI systems are providing a central hub for message transport, eliminating the need for custom transport code for each application. Up to now, however, EAI hub-and-spoke messaging has required a spaghetti of point-to-point analysis and mapping. The central information model works with EAI systems to provide coherence for data semantics, allowing *data* to become *information*.

ACKNOWLEDGEMENT

My thanks to Joram Borenstein for his helpful comments on this article.

AUTHOR BIOGRAPHY

Joshua Fox serves as Software Architect at Unicorn Solutions, working on the Unicorn Coherence™ platform for unification of information in the enterprise. Unicorn is a member of the W3C Web Ontology Working Group, which is standardizing ontology languages for the Web. Fox's previous experience includes the design and development of large-scale distributed Internet systems. He has earned a B.A. *summa cum laude* in Mathematics from Brandeis University and a Ph.D. in Comparative Semitic Philology from Harvard University, where he has also served as Lecturer. He has published and lectured extensively in the field of software engineering. Voice: 1-866-286-4267, ext 115; E-Mail: joshua@unicorn.com; Website: www.unicorn.com