

## Chapter 143

# Theory and Management of Data Semantics

**Daniel W. Gillman**

*Bureau of Labor Statistics, USA*

**Frank Farance**

*Farance Inc, New York, NY, USA*

*Category: Technologies for Knowledge Management*

### INTRODUCTION

Almost every organization, public or private, for profit or non-profit, manages data in some way. Data is a major corporate resource. It is produced, analyzed, stored, and disseminated. And, it is poorly documented.

Descriptions of data are essential for their proper understanding and use by people inside and outside the organization. For instance, systems for disseminating data on the Internet require these descriptions (Census Bureau, n.d.). Either

inside or outside the organization, functions of the system support finding the right data for a study, understanding data from a particular source, and comparing data across sources or time (Gillman, Appel, & LaPlant, 1996).

Descriptions of data and other resources are *metadata* (Gillman, 2003). Metadata are part of the corporate memory for the organization. Preserving corporate memory is one of the basic features of knowledge management (King, Marks, & McCoy, 2002). Metadata include the meaning, or semantics, of the data. In some countries, e.g., the U.S., a large percentage of the population is reaching retirement age. As a result, recording the memories of these workers, including the meaning of data, is increasingly important. Preserving metadata is crucial for understanding data years

DOI: 10.4018/978-1-59904-931-1.ch144

after the data were created (Gillman, Appel, & LaPlant, 1996).

Traditionally, the metadata for databases and files is developed individually, without reference to similar data in other sources. Even when metadata exist, they are often incomplete or incompatible across systems. As a result, the semantics of the data contained in these databases and files are poorly understood. In addition, the metadata often disappear after the data reach the end of the business life-cycle.

Techniques for documenting data are varied. There are CASE<sup>1</sup> tools such as Oracle Designer<sup>®</sup> (Oracle, n.d.) or Rational Rose<sup>®</sup> (IBM, n.d.). These tools produce models of data in databases (Ullman, 1982). The models provide some semantics for the data. For social science data sets, metadata is described in an XML<sup>2</sup> specification (ICPSR, n.d.). For geographic data sets, the U.S. Federal Geographic Data Committee developed a metadata framework, clearinghouse, and supporting software (FGDC, n.d.).

Metadata are data used to describe some objects. They are structured, semi-structured, or unstructured (Abiteboul, Buneman, & Suci, 2000), just as data are. Data are structured if one knows both the schema and datatype, semi-structured if one knows one of them, and unstructured otherwise. From the perspective of their content, documents are unstructured or semi structured data. Their schemas come from presentation frameworks such as HTML<sup>3</sup> (W3C, 1997) or word processor formats. Documents with the content marked up in XML (W3C, 2004) are semi-structured. When using the full datatyping capability of XML-Schema, the document is structured with respect to the content. However, the colloquial use of the term “document” begins to lose its meaning here.

In describing some resource, the content is more important than the presentation. The content contains the semantics associated with the resource. If the content is structured data, this increases the capability of performing complex

queries on it. Retrieving unstructured documents using search engine technology is not as precise.

It turns out there are structured ways to represent the semantics of data. Ontologies (Sowa, 2000) are the newest technique. Traditional database (or registry) models are examples of ontologies. This article describes the constituents of the semantics of data and a technique to manage them using a metadata registry. The process of registration – an approach to control the identification, provenance, and quality of the content – is also described and its benefits discussed.

## **SEMANTICS OF DATA**

### **Terminology**

The ancient Greek philosophers began the study of terminology and concept formation in language (Wedberg, 1982), and they discovered a useful relationship between term (word or phrase), concept, referent (object), and definition, that is illustrated in Figure 1 (CEN, 1995).

Figure 1 shows that concepts, terms, referents, and definitions are related but separate constructs. Each plays a role in our understanding (i.e., the semantics of) data.

An important observation is that concepts are human constructions (Lakoff, 2002). No matter how well we define a concept, a complete description is often impossible. Identifying the relevant characteristics is culturally dependent. So, some objects in the extension of a concept fit the characteristics better than others.

### **Framework for Terminology**

To begin, we describe some useful constructs from the theory of terminology. These come from several sources (Sager, 1990; Sager, 2000; ISO, 1999; ISO, 2000). We will use these constructs to describe the semantics of data. Some of the definitions have been slightly modified by the authors.

10 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

[www.igi-global.com/chapter/theory-management-data-semantics/49094](http://www.igi-global.com/chapter/theory-management-data-semantics/49094)

## Related Content

---

### Knowledge Transfer and Team Performance in Distributed Organizations

Ngoma Sylvestre Ngoma and Mary Lind (2015). *International Journal of Knowledge-Based Organizations* (pp. 58-80).

[www.irma-international.org/article/knowledge-transfer-and-team-performance-in-distributed-organizations/125585](http://www.irma-international.org/article/knowledge-transfer-and-team-performance-in-distributed-organizations/125585)

### Integrated Analysis and Design of Knowledge Systems and Processes

Mark E. Nissen, Magdi Kamel and Kishore Sengupta (2000). *Knowledge Management and Virtual Organizations* (pp. 214-244).

[www.irma-international.org/chapter/integrated-analysis-design-knowledge-systems/54262](http://www.irma-international.org/chapter/integrated-analysis-design-knowledge-systems/54262)

### Impact of Trust on Communication in Global Virtual Teams

Päivi Lohikoski, Jaakko Kujala, Harri Haapasalo, Kirsi Aaltonen and Leena Ala-Mursula (2016). *International Journal of Knowledge-Based Organizations* (pp. 1-19).

[www.irma-international.org/article/impact-of-trust-on-communication-in-global-virtual-teams/143217](http://www.irma-international.org/article/impact-of-trust-on-communication-in-global-virtual-teams/143217)

### Proposal of Indicators for Intellectual Capital in Higher Education

Edgar Oliver Cardoso Espinosa (2019). *The Formation of Intellectual Capital and Its Ability to Transform Higher Education Institutions and the Knowledge Society* (pp. 232-246).

[www.irma-international.org/chapter/proposal-of-indicators-for-intellectual-capital-in-higher-education/231064](http://www.irma-international.org/chapter/proposal-of-indicators-for-intellectual-capital-in-higher-education/231064)

### Co-Worker Support and Communities of Practice: Mediating Role of Employee Personal Interaction

Anjali Dutta and Santosh Rangnekar (2022). *International Journal of Knowledge Management* (pp. 1-17).

[www.irma-international.org/article/co-worker-support-and-communities-of-practice/297607](http://www.irma-international.org/article/co-worker-support-and-communities-of-practice/297607)