

DRAFT: The Basic Representation Model for Digital Preservation and Information Organization

Karen M. Wickett

Contents

1	Introduction	2
2	Background	2
2.1	General Theories of information	2
2.2	Related models	4
2.2.1	FRBR and related models	4
2.2.2	OASIS Reference Model	4
2.2.3	CIDOC-CRM	4
2.3	Data Conservancy Data Concepts Group	5
3	The Basic Representation Model	5
3.1	BRM Entities	5
3.2	BRM Relationships	6
3.3	Interpretive Frames	6
4	Discussion	7
4.1	Levels of representation	7
4.2	Application Areas	7
4.2.1	Scientific data management and preservation	7
4.2.2	Bibliographic data and Text encoding management	8
5	Future work	8
6	Conclusion	8

Abstract

Information models to handle preservation and information organization, such as the Functional Requirements for Bibliographic Records, the OAIS Reference Model, and the CIDOC Conceptual Reference Model have faced major stumbling blocks in their adoption, due in part to the difficulty in applying them consistently to digital objects. This paper describes and analyzes those challenges, and presents The Basic Representation Model, a general framework to account for the multiple and varied levels of representation that occur in the creation, management, and preservation of digital objects. This model was developed in the context of a project focused on the preservation of scientific data, and has the additional potential to support fine-grained information modeling in bibliographic, scientific, and cultural heritage domains. In particular, the model can support clear and explicit connections between information organization, information management, and digital preservation.

1 Introduction

[Section goals: The potentials for connecting information organization, information management, and digital preservation. My goal with the paper is to provide a general model of the objects and relationships that are fundamental to the representation and transmission of information. My goal with this section is to motivate the need for such models. I think what is missing in many models is a flexibility to handle the multi-level realization and embodiment relationships that typically occur in digital systems.]

2 Background

My goal in this section is to sketch the intellectual history of the Basic Representation Model and review and critique the related models that were the primary inspiration for BRM.

2.1 General Theories of information

The development of conceptual models to explain the entities and relationships involved in the recording and transmission of information is an important activity for the field of library and information science. These models inform system design, guide research, and provide insight into the nature of information and information-bearing objects. The Basic Representation Model bridges the gap between general models of the what information is, and the encoding and transmission of information into documents and digital objects.

At their most general, theories of information are oriented around communication and the transmission of messages that carry content. For example, the Mathematical Theory of Communication (Shannon and Weaver, 1964) views information content as quantifiable in terms of the reduction in uncertainty given by the reception of a signal. While this model encapsulates an engineering approach to building robust channels for signal transmission and does not accord well with the realities of human communication (Ma, 2012), it highlights the fact that information consists in some sense of content.

The centrality of information to human communication has driven many philosophical explorations of the topic. The theories developed by Barwise and Perry (1983), Dretske (1981), and Devlin (1995) characterize meaning and information as highly contextualized phenomena. Information-bearing messages have meaningful content only and because of human action - in expression of content and in the interpretation of signals. Floridi brings together the philosophical and engineering approaches in his General Definition of Information, which treats information as consisting of some quantifiable amount of data, which is well-formed by some syntax and meaningful under some interpretation (Floridi, 2013). Floridi's and other philosophi-

cal accounts are explanatory for understanding the basic nature of information and its fundamental role in human communication, but do not connect to the management of information as approached in library and information studies.

Bates’s “Fundamental Forms of Information” takes a similarly generalist approach and presents a wide-ranging analysis of information in its many forms. Bates advocates for a suite of general definitions for information, with her most general account defining information as a “pattern of organization of matter and energy” (Bates, 2006). The general account is divided into several classes of information, with the most relevant for the work discussed here being her definition of Recorded Information as “Communicatory or memorial information preserved in a durable medium.” The management of records and the information they contain is a core activity for library and information science, and The Basic Representation Model described below is concerned with information in this sense.

In “Information as Thing”, Michael Buckland presents a view of the information field that analyzes the work, outputs and perspectives on information in terms of “information as thing”, “information as process” and “information as knowledge”. Buckland points out that “Information-as-thing is of special interest in the study of information systems. It is with information in this sense that information systems deal directly.” The conceptual models of information-bearing objects (discussed below) that have been developed to promote interoperability and system-building since them all confirm Buckland’s stance on information (as thing) by identifying entities and relationships pertaining to the expression, encoding and storage of information in some concrete form.

One of the primary entities identified by Buckland as a focal point for the analysis of information systems is the document. Although the precise nature of documents has also been argued to be contextual (cite Buckland 1997), Elaine Svenonius grounds her analysis of information organization systems in the notion of the document, defined as “an information-bearing message in recorded form”. Svenonius then goes on to define and analyze information organization systems in terms of set operations on documents as determined by various features. Svenonius’ account is explanatory of the purposes and activities that are fundamental to information organization, but does not analyze documents as bearing information in recorded form.

[TRANSITION]

2.2 Related models

2.2.1 FRBR and related models

The Functional Requirements for Bibliographic Records (FRBR) was produced by IFLA in 1998 and has been a focal point in the evolution of bibliographic cataloging since then (IFLA, 2009). Many folks have thought about how to apply it to digital objects in general (Renear Dubin ASIST 2007)(Renear and Dubin, 2007), video games and virtual worlds (McDonough et al 2010)(McDonough et al., 2010), scientific data (Hourcle 2008, DCDC ASIST 2011,2012)(Hourcl , 2008) (Sacchi et al., 2011b).

FRBR Group One entities and relationships. Since FRBR presents these in pretty agnostic terms, it tempting to apply the relationships outside the scope of monographic publishing. But it's not clear that FRBR even works that well for tracking the entities and relationships involved in, for example, the publication of XML documents, or e-books. Certainly problems arise when you try to apply it to scientific data, see (Wickett et al., 2012).

2.2.2 OASIS Reference Model

So, OASIS is really important and influential for digital preservation systems (CCSDS, 2002). Everybody be making their SIPs and their DIPs. And it seems to work really well for the one level of representation it works at, which is files. But then it also gives a few models of information objects and content that are not easy to understand and don't work as general models. As Simone put it, One Thing is Missing or Two things are confused (Sacchi et al., 2011a).

An incomplete list of important concepts to explain and critique from OASIS: information object, data object, representation information, designated community.

2.2.3 CIDOC-CRM

CIDOC-CRM is an ontology for descriptions of cultural heritage objects. I'm including it here because it aims to give a high level ontology that museum descriptions can be mapped to to enable interoperability. I see a similar potential role for BRM but with scientific data. CIDOC is very oriented around cultural heritage and museum object description practices.

This makes me wonder though, about how BRM maps to CRM. Also can BRM serve as a sort intermediate layer that would enable mapping of scientific data to various scientific ontologies, like SWEET, etc.

2.3 Data Conservancy Data Concepts Group

The Data Conservancy Data Concepts (DCDC) group was a research group hosted at the Center for Informatics Research in Science and Scholarship at the University of Illinois at Urbana-Champaign, as part of the NSF Data Conservancy project. The aims of the research group were to clarify key concepts in the encoding and representation of scientific research data. As part of this work, the group worked to identify appropriate models enable description and preservation of digital datasets that were the result of scientific data collection, observation, analysis, and aggregation processes.

These mapping efforts focused on popular models from library and information science and resulted in a series of papers (Sacchi et al. (2011b), Dubin et al. (2011), Wickett et al. (2012)). The group proposed the Systematic Assertion Model (SAM), a conceptual model for scientific data that positioned the content of scientific data as propositional content warranted by an observational or computational event, and data as the symbolic representation of that data content.

The Basic Representation Model provides the representational backbone for the Systematic Assertion Model. SAM is an event-based model that accounts for the agents and events that essential to the origination of scientific data. BRM provides the entities and relationships that those agents and events act on or produce. The popular models that the group reviewed in Sacchi et al. (2011b) all had big problems for us, which I'll tell you about in the next subsection. They didn't have enough fine-grained divisions, or they seemed to conflate important objects, relationships, or activities.

The entities and relationships in BRM were first named and published on by Dave Dubin as part of Preservation Model 1, which was a result of the EchoDep project (Sandore and Unsworth, 2010).

3 The Basic Representation Model

3.1 BRM Entities

Propositional content. Allen Renear likes to sometimes define propositions as the bearers of truth values. Propositional content is the kind of thing (and the only kind of thing) that can be true or false. But the concern here is not with whether those propositions are true or false, or whether they pertain to our actual world. Just that they are purely content.

Symbol Structure. These are arrangements of elements into discernable patterns.

Patterned matter and energy. We are living in a material world, and to encounter information we

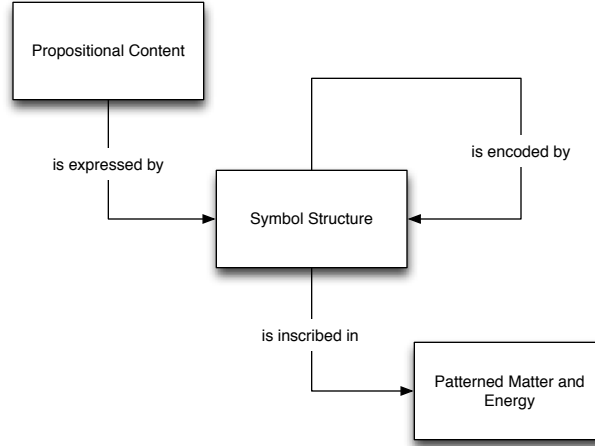


Figure 1: The Basic Representation Model

must encounter some material objects that record it in a way we can interpret.

3.2 BRM Relationships

Is Expressed By. This seems to be lifted from FRBR. This relationship can stand between propositional content and symbol structures. The symbol structure that stands in this relationship is a *primary symbol structure*. The distinction between primary and non-primary symbol structures was important for identifying data content.

Is Encoded By. this is a relationship between symbol structures. It can describe encoding relationships in digital systems. It is the possibility of recording information about relationships between symbol structures in the expression and encoding of information that is particularly useful for modeling information in digital computing systems. These layers of symbol structures in encoding relationships are key to the recording and storage of information in computational systems.

Is Inscribed In. This is the relationship between a symbol structure that expresses content and a concrete arrangement of matter or energy.

3.3 Interpretive Frames

Text from Dubin et al. (2011):

But the digital data resources that concern us are encoded symbol structures that express data content. Our problem is the contingent nature of this connection: data express their conceptual

content not simply in virtue of their arrangement and structure, but always with reference to what we call *interpretive frames*. These are abstractions that frame the interpretive context for symbolic expressions.

At the risk of understating their complexity, one can think of interpretive frames as functions or mappings between structural propositions at different expressive levels, or from structural propositions to conceptual propositions. Examples of interpretive frames include the grammatical rules expressed by an XML Schema, coded character sets such as ACSCII, the convention of writing numbers as strings of Arabic numerals with ten as the implied numerical base, the Hierarchical Data Format standard, and all dialects of the English language as they are spoken today. Interpretive frames also include any systematic expressive choices that may be local to a particular digital resource, such as a correspondence between successive rows of a spreadsheet and the order of transactions in a scientific experiment.

4 Discussion

4.1 Levels of representation

Information models that assume a single level (or a fixed set of levels) of representation in digital objects are always going to suck to apply. I'm not sure this needs to be a separate section, really, or whether it might just be most effective to critique the representation of digital objects by the various models in the background section.

4.2 Application Areas

4.2.1 Scientific data management and preservation

Ideas for application scenarios to discuss:

- instrument data
- astronomical data
- biodiversity records

4.2.2 Bibliographic data and Text encoding management

- Jacob and Dave’s case study from Balisage 2018.
- MARC records into XML into linked open data

5 Future work

- Integration with provenance models

6 Conclusion

References

J. Barwise and J. Perry. *Situations and Attitudes*. MIT Press, 1983.

Marcia J. Bates. Fundamental forms of information. *Journal of the American Society for Information Science and Technology*, 57(8):1033–1045, 2006. doi: 10.1002/asi.20369. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.20369>.

CCSDS. Reference model for an open archival information system (OAIS). Technical report, CCSDS 650.0-B-1, Blue Book, 2002.

Keith Devlin. *Logic and information*. Cambridge University Press, 1995.

F.I. Dretske. *Knowledge and the Flow of Information*. Blackwell, 1981.

David Dubin, Karen M. Wickett, and Simone Sacchi. Content, format, and interpretation. *Proceedings of Balisage: The Markup Conference 2011*, 2011. doi: 10.4242/balisagevol7.dubin01. URL <http://dx.doi.org/10.4242/BalisageVol7.Dubin01>.

Luciano Floridi. *The philosophy of information*. OUP Oxford, 2013.

Joseph A. Hourclé. FRBR applied to scientific data. In *Proceedings of the American Society for Information Science and Technology*, volume 45, pages 1–4, 2008. doi: 10.1002/meet.2008.14504503102. URL <http://dx.doi.org/10.1002/meet.2008.14504503102>.

- IFLA. Functional requirements for bibliographic records: Final report. Technical report, International Federation of Library Associations and Institutions, 2009. URL <http://www.ifla.org/files/cataloguing/frbr/frbr2008.pdf>.
- Lai Ma. Meanings of information: The assumptions and research consequences of three foundational theories. *Journal of the American Society for Information Science and Technology*, 63(4):716–723, 2012.
- Jerome McDonough, Matthew Kirschenbaum, Doug Reside, Neil Fraistat, and Dennis Jerz. Twisty little passages almost all alike: Applying the frbr model to a classic computer game. *Digital Humanities Quarterly*, 4(2):1869–1883, 2010.
- A. H Renear and D. Dubin. Three of the four FRBR group 1 entity types are roles, not types. *Proceedings of the American Society for Information Science and Technology*, 44(1):1–19, 2007.
- Simone Sacchi, Karen M. Wickett, Allen H. Renear, and David Dubin. One thing is missing or two things are confused: An analysis of OAIS Representation Information. *Poster presented at the Seventh International Digital Curation Conference*, 2011a.
- Simone Sacchi, Karen M. Wickett, Allen H. Renear, and David S. Dubin. A framework for applying the concept of significant properties to datasets. In *Proceedings of the American Society for Information Science and Technology*, New Orleans, LA, 2011b.
- Beth Sandore and John Unsworth. *ECHO DEPository — Phase 2: 2008–2010 Final Report of Project Activities*, section 4.2, pages 30–37. University of Illinois at Urbana-Champaign, June 2010.
- C. E. Shannon and W. Weaver. *The Mathematical Theory of Communication*. The University of Illinois Press, 1964.
- Karen M Wickett, Simone Sacchi, David Dubin, and Allen H Renear. Identifying content and levels of representation in scientific data. *Proceedings of the American Society for Information Science and Technology*, 49(1):1–10, 2012.