

"UNSTRUCTURED DATA" PRACTICES IN POLAR INSTITUTIONS AND NETWORKS: A CASE STUDY WITH THE ARCTIC OPTIONS PROJECT

Paul Arthur Berkman^{1,2*}

**¹ Bren School of Environmental Science and Management / Marine Science Institute, University of California, Santa Barbara, CA 93106, USA*

Email: berkman@bren.ucsb.edu

² DigIn (Digital Integration Technology Limited), 6 Caxton House, Broad Street, Great Cambourne, Cambridge, CB23 6JN, UK

Email: paul.berkman@digin.co.uk

ABSTRACT

Arctic Options: Holistic Integration for Arctic Coastal-Marine Sustainability is a new three-year research project to assess future infrastructure associated with the Arctic Ocean regarding: (1) natural and living environment; (2) built environment; (3) natural resource development; and (4) governance. For the assessments, Arctic Options will generate objective relational schema from numeric data as well as textual data. This paper will focus on the "long tail of smaller, heterogeneous and often unstructured datasets" that "usually receive minimal data management consideration," as observed in the 2013 Communiqué from the International Forum on Polar Data Activities in Global Data Systems.

Keywords: Big Data, unstructured data, relational schema, infrastructure, integration

1 Introduction

Interests are awakening globally to take advantage of extensive energy, shipping, fishing, and tourism opportunities associated with diminishing sea ice in the Arctic Ocean (Berkman and Vylegzhanin, 2013). Recognizing these diverse interests, there is urgency to develop infrastructure so the commercial activities can proceed in a sustainable manner, balancing:

- National interests and common interests;
- Environmental protection, social equity and economic prosperity; and
- Needs of present and future generations.

Infrastructure in the Arctic Ocean will include port facilities, sea lanes, emergency response assets, communication systems and observing networks as well as regulatory and policy systems. Cross-cutting all aspects of the infrastructure will be information management and knowledge discovery using disparate data. The *Arctic Options: Holistic Integration for Arctic Coastal-Marine Sustainability* project, which is being funded by the National Science Foundation in the United States and Centre Nationale de la Recherche Scientifique in France from 2013-2016, provides a case study for international, interdisciplinary and inclusive data practices.

As part of the *Arctic Science, Engineering and Education for Sustainability* programme (ArcSEES 2012), *Arctic Options* will consider data for:

1. Natural and living environment;
2. Built environment;
3. Natural resource development; and
4. Governance.

To enhance its cost-effectiveness, *Arctic Options* also has established links to the *Study of Environmental Arctic Change* (SEARCH 2013) and *Arctic Climate Change, Economy and Society* (ACCESS 2013) projects that are supported extensively within the United States and Europe, respectively.

These data for Arctic Ocean infrastructure will be generated from sensor and transactional systems from observing networks (AOS 2013) as well as experiments in numeric formats that are "structured" (i.e., managed) with databases, which can be analyzed statistically and graphically with various relational approaches. Geographic Information Systems (GIS) will be particularly powerful for marine spatial planning, ecosystem-based management and integrated ocean management in the Arctic Ocean (Håkon Hoel, 2010; Ehler, 2011; Clement et al., 2013; PAME, 2013).

In addition, the data in the *Arctic Options* project will involve digital resources in natural language formats (e.g., papers, reports and agreements), which commonly are considered to be "unstructured" (i.e., unmanaged) because they "*cannot be decomposed into standard components*" or relational schema (Oracle, 2002). The unstructured data will be aggregated from diverse institutions that have Arctic remits (Berkman and Vylezhanin, 2013), such as the Arctic Council (2013).

Within the popular framework of 'Big Data' (Lohr, 2012), structured and unstructured data together reflect the full complement of digital information that we produce as a global society, with the volume of unstructured data accounting for upwards of 85% of the information and growing twice as fast as structured data (Fig. 1). In this paper, innovations with unstructured data will expand on earlier developments through the National Science Digital Library (NSDL, 2013) and International Research on Permanent Authentic Records in Electronic Systems project (InterPARES, 2013), as summarized in a publication through the Committee on Data for Science and Technology (Berkman et al., 2006).

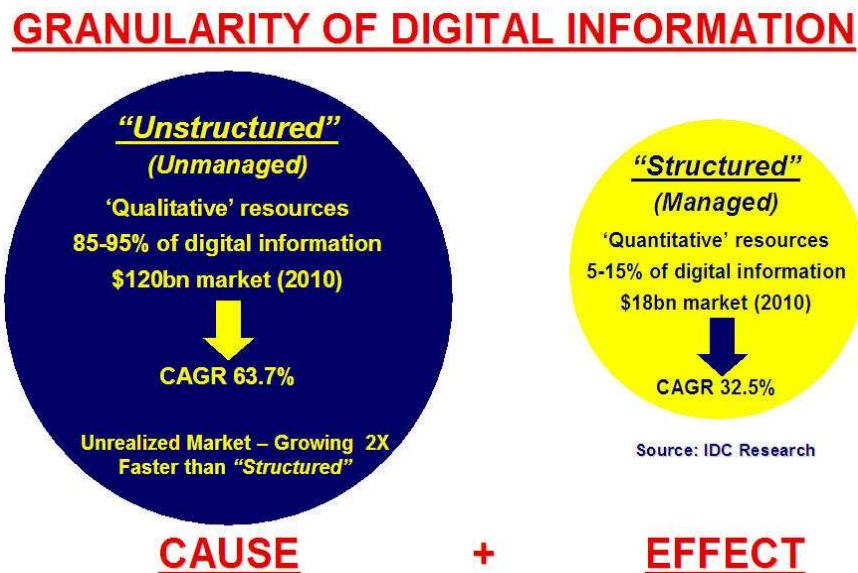


Figure 1. 'Big Data' defined in terms of "structured" and "unstructured" data, both of which relate to granularity of the information resources. Compound annual growth rates (CAGR) and estimated market sizes are from IDC (2010).

This paper also will focus on manipulation of unstructured data in view of the following observation in the *Communiqué* from the *International Forum on Polar Data Activities in Global Data Systems* (Polar Data Forum, 2013):

It is the long tail of smaller, heterogeneous and often unstructured datasets (those without metadata, mark-up and not in databases) that receive the least data management attention by scientific repositories. However, utilizing the inherent structure of any digital resources provides an objective framework to discover their relationships in a manner that complements existing content and context management solutions.

In particular, this paper is intended to provoke discussion about applying the inherent structure of digital information, which is a fundamental opportunity with digital information and distinction compared to all hardcopy resources.

2 Elements of Meaning

Each era of global communication, from stone to digital (Berkman 2008), has been accompanied by a threshold increase in human capacity to transport information. Similarly, each new communication medium has significantly increased our capacity to produce information, as indicated by the relative volumes of information that emerged. Moreover, the ability to integrate information has increased over time with tablets, folios, books, and now websites. In contrast, the most resilient medium was stone with petroglyphs and pictographs that have stood the test of time through rain, snow, wind, and even fire. Subsequent media have been much more fragile. In fact, the digital medium has been like a black hole where most of the information produced has been lost because of limited preservation strategies and rapid obsolescence of storage devices.

Looking backward through time, information in our civilization has been managed largely through libraries and archives. While similar in their needs to facilitate information access and preservation, these two architectures possess fundamental differences. Archives manage information based on the *context* of records linked to specific activities and transactions, like the housing authority that records the title of your home. Libraries largely manage information based on the *content* of the information resources, as with the subject categories in the Dewey Decimal System (OCLC, 2013).

Beyond content and context, the third element of information to establish meaning is its *structure*. For example, when a message is encrypted (i.e., the structure is altered) it still has content and context, but no meaning absent the key to unlock the encryption. Alternatively, if the names or dates and places are removed from an information resource, it still has context and structure, but limited meaning without the salient facts. Similarly, meaning will be compromised by removing the context that can be used to authenticate an information resource or establish its provenance.

All information must have content, context and structure to create meaning (Fig. 2). However, with the digital medium it is possible to utilize the *content* and *context* as well as *structural* patterns to manage sets, subsets, and supersets of information resources. **The capacity to utilize the inherent structure of digital resources is the distinguishing feature of digital information compared to all of its hardcopy predecessors.**

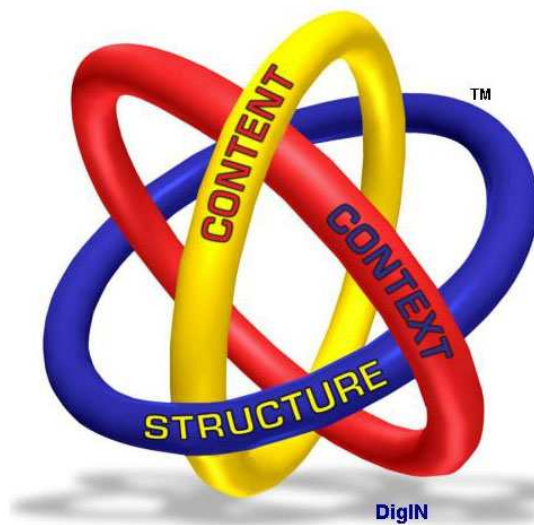


Figure 2. Borromean ring, illustrating the interconnected core elements of all information - both hardcopy and digital - that together create meaning (revised from Berkman, 2008).

3 Digital Information Architectures

In general, unstructured data are managed with metadata, markup or databases. However, for text resources specifically, the digital resource itself contains the information content that would be summarized by metadata. Consequently, metadata incompletely and subjectively characterizes digital text resources for the purposes of information access, as card catalogues did with hardcopy resources (OCLC 2013).

Nonetheless, ubiquitous use of metadata, which originated with card catalogues for hardcopy libraries (Dublin Core 2003), has become a *de facto* approach for digital information management around the world with diverse

"standards" through the International Organization for Standardization (e.g., ISO 2013). These standards vary by country and require extensive effort in terms of personnel expertise, time and cost to implement.

Although not properly quantified, back-of-the-envelope calculations further suggest that metadata production may account for more than 10% of the global expenditure on information and communications technology. Most importantly, the production of metadata does not scale with increasing granularity (Berkman et al. 2006), which largely explains the growing discrepancy between the volumes of unstructured and structured digital information (Fig. 1).

Data that is unmanaged with current technologies raises the question about strategies for information management and knowledge discovery with text resources. Given the global diversity of information technology companies, Table 1 was created to compare attributes and functions for information management and knowledge discovery (Table 2) across generalized solution suites that apply the content and context as well as the structure of digital text resources.

Table 1 Capacity to Utilize Various Attributes and Functions (Yes or No) in Relation to Generalized Solution Suites for Digital Information Management and Knowledge Discovery

| Attributes and Functions (see Table 2) | Interconnected Core Elements of Meaning (Fig. 2) and Underlying Solution Suites for Digital Information Management and Knowledge Discovery | | | |
|---|--|--------------------------|--|---------------------|
| | Content and Context | | | Structure |
| | Search Engines | Databases / Spreadsheets | Metadata / Ontologies / Semantic Indexes | Granularity Engines |
| Language Independent | N | Y | N | Y |
| File-Type Independent | Y | N | N | Y |
| Scale Independent | Y | N | N | Y |
| Classification Independent | N | Y | N | Y |
| Ranking Independent | N | Y | Y | Y |
| Markup Independent | Y | Y | N | Y |
| Metadata Independent | N | Y | N | Y |
| Re-Purpose Metadata Tags | N | N | N | Y |
| Tabular Manipulation | N | Y | N | Y |
| Result Lists | Y | Y | Y | Y |
| Relational Displays | N | Y | Y | Y |
| Relational Analytics | N | Y | Y | Y |
| Preserves Authentic Record | Y | N | N | Y |
| Content-Driven | Y | N | Y | Y |
| Context-Driven | N | Y | Y | Y |
| Structure-Driven | N | Y | N | Y |
| User-Defined Rules | N | Y | Y | Y |
| Single-Level Inverted Indexing | Y | Y | Y | Y |
| Multi-Level Inverted Indexing | N | N | N | Y |
| 2 ^N Permutations | N | N | N | Y |
| Result-Set Certainty | N | N | N | Y |
| Automated Granularity | N | N | N | Y |

Table 2 Descriptions of Attributes and Functions in Table 1

| Attribute or Function | Description |
|---|--|
| Language Independent | Operations not limited by symbolisms, such as different alphabets |
| File-Type Independent | Operations not limited by types of digital resources, recognizing that all resources (e.g., text, images, genomes, sensor data streams, music) have structural patterns of embedded organization that can be manipulated automatically |
| Scale Independent | Operations not limited by size of resource set, which can be applied by an individual, small business, large corporation or government |
| Classification Independent | Operations not limited by subjective categories that are defined by individuals or programmed algorithms |
| Ranking Independent | Lists generated without programmer algorithms that arbitrarily rank relevant search results |
| Markup Independent | Navigation not limited by markup tagging, which is subjective |
| Metadata Independent | Access not limited by metadata schema because all symbols in a resource set are indexed and searchable at all levels of embedded granularity |
| Re-Purpose Metadata Tags | Applies existing metadata fields associated with a digital resource to manipulate the embedded levels of organization in expandable-collapsible hierarchies |
| Tabular Manipulation | Capability to combine rows, columns and cells from multiple tables |
| Result Lists | Capability to generate linear displays of search results |
| Relational Displays | Capability to generate integrated displays of search results |
| Relational Analytics; | Capability to generate statistical displays of integrated search results |
| Preserves Authentic Record | Capability to preserve authentic digital resources while at the same time providing the ability to search or integrate the digital resources |
| Content-Driven | Based on information classification schema that are subjectively defined by individuals or programmed algorithms |
| Context-Driven | Based on information provenances that are subjectively defined by individuals or programmed algorithms |
| Structure-Driven | Based on information boundaries and patterns that can be manipulated within and between embedded levels of organization in digital resources |
| User-Defined Rules | Contrasted with programmer-constrained rules |
| Single-Level Inverted Indexing | One to many relationships among digital objects that is used to generate lists |
| Multi-Level (Dynamic) Inverted Indexing | Many to many relationships within and between digital objects that is used to generate expandable-collapsible hierarchies |
| 2 ^N Permutations | Comprehensive capacity to integrate 'N' digital objects and discover all possible combinations of granules within and between digital resources |
| Result-Set Certainty | Objective and comprehensive results in contrast to probabilistic solutions, which have inherent uncertainties |
| Automated Granularity | Generate subsets of embedded parent-child relationships down to finite elements at the lowest levels of granularity within and between digital resources |

Among the solutions in Table 1 - search engines, databases and spreadsheets are well known and need no elaboration. Similarly, ontologies and semantic indices are widely used (Berners-Lee et al. 2001). The concept of a "granularity engine," however, is being introduced herein as a fundamental solution based on the inherent structure of digital resources. In particular, granularity engines can leverage the inherent structure of digital resources to achieve functionalities beyond what is possible with solutions derived from their content or context (ie., described by the full complement of attributes and functions with "Y" in Table 1).

Most notably, granularity engines can objectively deliver 2^N relationships among N digital objects, overcoming a significant shortcoming with subjective content or context solutions that limit the range of relationships that can be discovered. Consider two digital objects where the four possible permutations include one, the other, both or neither. If there are just one hundred digital objects, which is a small number, the number of possible permutations (2^{100}) effectively would be a googol and we have the challenge to discover relationships across thousands and millions of digital objects.

4. Inherent Structure and Granularity Engines

Effectively, we all have infinite and instantaneous access to digital data on our computers or networks and over the internet. With text resources, searching through repositories merely lists items that contain the search query. Lists generally are ranked in an arbitrary manner, commonly in terms of assumed relevance. From lists of possibly relevant results, digital resources can be selected and it is then up to the user to hunt sequentially for the search term through each resource. If the user wants to identify content-in-context relationships within and between the resources (e.g., relevant sentences within chapters or resources within years), it is then necessary to: (a) cut the relevant pieces out of each relevant resource; (b) paste them into a new folder; and then (c) organize all of the cut-and-paste pieces. This a-b-c process to establish relationships within and between digital resources is tedious, time consuming and subjective.

However, text resources have structure that is defined by the grammar rules of the language. For example – read left-right and top-bottom in English – books have chapters that have pages with paragraphs embedded with sentences composed of words that each contain letters. Through each book, various headings and forms of punctuation (e.g., period, exclamation point or question mark at the end of a sentence) define boundaries that can be used to disaggregate embedded levels of granularity, like peeling an onion.. Unlike hardcopy resources, repeating structural patterns in digital resources can be set as rules (Berkman et al. 2006) to run a granularity engine that generates and indexes discrete granules (e.g., sentences, paragraphs, pages or chapters). Subsequently, these granules can be integrated in parent-child contexts across a collection of digital resources for any search query. Such management, discovery and analysis of digital text documents with a granularity engine will complement the GIS manipulations of numeric data layers in the *Arctic Options* project.

A classic form of an unstructured resource is a PDF (portal document format) file, which has advantages of being interoperable across diverse operating systems and file formats as well as serving as an archival standard (ISO 2009). Utilizing a familiar collection to illustrate the application of a granularity engine, 53 PDF files of books written by Charles Dickens from 1836-1880 were automatically decomposed into 571,386 granules that represent all sentences, paragraphs, pages and chapters within these books and years (Fig. 3).



Next Generation - Simple Discovery Dickens Christmas Present



best of times Exact Match

Digital Zoom™ Year Resource Section Page Paragraph Sentence

15 Years - 16 Resources - 15 Sections - 18 Pages - 18 Paragraphs - 18 Sentences

- Publication Years
 - 1836
 - The Pickwick Papers.pdf
 - 1837
 - Oliver Twist.pdf
 - 1838
 - Nicholas Nickleby.pdf
 - 1840
 - 1842
 - 1843
 - 1844
 - 1846
 - 1856
 - 1859
 - A Tale Of Two Cities.pdf
 - Chapter I - The Period
 - 2
 - Para 1
 - It was the **best of times** it was the worst of times it was the age of wisdom it was the age of foolishness it was the epoch of belief it was the epoch of incredulity it was season of Light it was the season of Darkness it was the spring of hope it was the winter of despair we had everything before us we had nothing before us we were all going direct to Heaven we were all going direct the other way - in short the period was so far like the present period that some of its noisiest authorities insisted on its being received for good or for evil in the superlative degree of comparison only.
 - Chapter XV - Knitting
 - 154
 - Para 1
 - Monsieur Defarge sold a very thin wine at the **best of times** but it would seem to have been an unusually thin wine that he sold at this time.
 - 1860
 - 1861
 - 1864
 - 1869
 - 1880

Figure 3. Granularity engine implementation with 53 53 PDF files of books written by Charles Dickens from 1836-1880, which were automatically decomposed by PDF KnoHow™ (DigIn 2013) into 571,386 granules that represent all sentences, paragraphs, pages and chapters within these books and years. The exact match search for "best of times" reveals occurrence of this famous phrase in 15 years, 16 books, 15 chapters, 18 pages, 18 paragraphs and 18 sentences that can be expanded and collapsed with the Digital Zoom™. The relevant granules can be aggregated and further analysed to quantify parent-child frequencies comprehensively at all granularity levels in the expandable-collapsible hierarchy.

In the *Arctic Options* project, complementing the geospatial integration of data layers with GIS, granularity-engine applications will enable users to objectively zoom in and out of content layers across collections of digital text resources for any Boolean search query. As an example, to discover relationships across the entire Dickens collection, searching for "best of times" in Tale of Two Cities surprisingly reveals that this famous phrase was repeated within Dickens' books throughout his career (Fig. 3). Such surprises, which are at the heart of discovery, can be generated by granularity engines based on the inherent structure of digital resources without metadata, markup or databases (Table 1).

5 Conclusions

There is no such thing as "unstructured" data, because all data must have structure as well as content and context to have meaning (Fig. 2). Moreover, granularity-engine applications with PDF files (Fig. 3) falsify long-standing definitions for unstructured data (e.g., Oracle, 2002) because they can be automatically "*decomposed into standard components*" as well as relational schema without metadata, markup or databases. The unique advantage of digital resources over hardcopy resources is the opportunity to utilize the inherent structure of the resources as well as their content and context for the purposes of information management and knowledge discovery.

6 Acknowledgements

This paper was presented as invited keynote at the International Forum on Polar Data Activities in Global Data Systems (Polar Data Forum, 2013) and has been prepared with support from *Arctic Options: Holistic*

Integration for Arctic Coastal-Marine Sustainability project, supported by the National Science Foundation (NSF-PLR No. -1263819).

7 References

ACCESS 2013. Arctic Climate Change, Economy and Society. European Commission (<http://www.access-eu.org>; accessed 15 December 2013)

Arctic Council. 2013. (<http://www.arctic-council.org>; accessed 28 January 2014).

ArcSEES 2012. *Arctic Science, Engineering and Education for Sustainability*. National Science Foundation, (<http://nsf.gov/pubs/2012/nsf12553/nsf12553.htm>; accessed: September 2013).

Berkman, P.A. 2008. Once in a hundred generations. In: Halbert, M. and Skinner, K. (eds.). *Strategies for Sustaining Digital Libraries*. Emory University, Atlanta., Georgia, USA. Pp. 11-21.

Berkman, P.A., & Vylegzhanin, A.N. (eds.). 2013. *Environmental Security in the Arctic Ocean*. Springer, Dordrecht, The Netherlands.

Berkman, P.A., Morgan, G.J., Moore, R., & Hamidzadeh, B. 2006. Automated Granularity to Integrate Digital Information: The “Antarctic Treaty Searchable Database” Case Study. *Data Science Journal* 5:84-99. (https://www.jstage.jst.go.jp/article/dsj/5/0/5_0_84/article. accessed: September 2013).

Berners-Lee, T., Hendler, J., & Lassila, O. 2001. The Semantic Web. *Scientific American*, May 2001, p. 29-37.

Clement, J.P., Bengston, J.L. and Kelly, B.P. 2013. *Managing for the Future in a Rapidly Changing Arctic. A Report to the President*. Interagency Working Group on Coordination of Domestic Energy Development and Permitting in Alaska. Washington, D.C., USA.

DigIn. 2013. Dickens Christmas Present. Digital Integration Technology Limited (<http://dickens.knohow.co>; accessed 15 December 2013).

Dublin Core. 2003. Dublin Core Metadata Initiative. (<http://dublincore.org/resources/faq/>; accessed 15 December 2013).

Ehler, C. 2011. Part II. Marine Spatial Planning in the Arctic: A first step toward ecosystem-based management. In: *The Shared Future. A Report of the Aspen Institute Commission on Arctic Climate Change*. The Aspen Institute, Washington, D.C. pp. 40-82. (<http://www.unesco-ioc-marinesp.be/publications>; accessed 15 December 2013).

Håkon Hoel, A. 2010. Integrated Oceans Management in the Arctic: Norway and Beyond. *Arctic Review on Law and Politics*, 1:186-206.

IDC. 2010. Digital Universe Study, A Digital Universe Decade – Are Your Ready? IDC Corporation.

InterPARES. 2013. International Research on Permanent Authentic Records Project (<http://interpares.org>; accessed 15 December 2013).

ISO. 2009. ISO 19005-1:2005. Document management -- Electronic document file format for long-term preservation -- Part 1: Use of PDF 1.4 (PDF/A-1). (http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=38920; accessed 29 January 2014).

ISO. 2013. ISO/IEC JTC1 SC32 WG2 Development/Maintenance: ISO/IEC 11179, Information Technology -- Metadata registries (MDR). (<http://metadata-standards.org/11179>; 15 December 2013).

Lohr, S. 2012. The Age of Big Data. *New York Times*, 11 February 2012.

NSDL 2013. National Science Digital Library (<http://www.nsdl.org>; accessed 15 December 2013).

OCLC. 2013. Dewey Decimal Classification (<http://www.oclc.org/dewey.en.html>; accessed 29 January 2014).

Oracle. 2002. Oracle9i Application Developer's Guide - Large Objects (LOBs). Release 2 (9.2). Oracle Corporation. (http://docs.oracle.com/cd/B10501_01/appdev.920/a96591/adl01int.htm; accessed 15 December 2013).

PAME. 2013. Arctic Ocean Review. Final Report. Phase II: 2011-2013. Arctic Council Working Group on the Protection of the Arctic Marine Ecosystem (PAME), Akureyri, Iceland.

Polar Data Forum. 2013. Communiqué from the *International Forum on Polar Data Activities in Global Data Systems*. International Council of Science (World Data System, Scientific Committee on Antarctic Research, International Arctic Science Committee, Committee for Data Science and Technology), 15-16 October 2013, Tokyo, Japan.

SEARCH 2013. Study of Environmental Arctic Change. Arctic Research Consortium of the United States (<http://www.arcus.org/search>; accessed 15 December 2013)