Chapter 38 Massive Digital Libraries (MDLs)

Andrew Philip Weiss California State University – Northridge, USA

ABSTRACT

Massive digital library (MDL) is a term coined to define a class of digital libraries gathering mass-digitized print books and monographs, which rival the size of brick-and-mortar libraries. Specific examples of MDLs, including Google Books, HathiTrust, DPLA, Internet Archive, et al., are presented. The issues raised by MDLs include mass-aggregation of digital content and the ability to maintain source-material accuracy and veracity; copyright, fair use, and the mass-digitization of materials not in the public domain; and disparities in the level of diversity, especially with regard to Spanish-language, Japanese-language, and Hawaii-Pacific materials. Finally, the impact of MDLs on Digital Humanities, especially with regard to the Google Books digital corpus and the Google Ngram Viewer, will be investigated.

INTRODUCTION

To provide a clearer framework for analyzing the growth of digital libraries, Weiss and James have proposed the term Massive Digital Libraries (MDLs), which is based on the size, scope and increasing scalability of digitized book collections. Such MDLs rival the size, breadth, and depth of a physical library's print holdings, and often reach a scale seen among library consortia collections. (Weiss and James, 2013a, 2013b, 2014, 2015; Weiss, 2016)

The root of the concept begins in late 2004 when Google made its "resounding announcement" to digitize millions of the world's books—including works still under copyright protection—and to place them *all* online. (Jeanneney, 2005) Jean-Noel Jeanneney, head of Bibliothèque nationale de France at the time, interpreted Google's planned project as a wake-up call for European countries. Failure to catch up to the American company, he argued, would result in significant problems for non-American organizations.

Twelve years on, it is hard to imagine that Google's desire to create an online digital library on such a large scale should have come as such a shock. Yet at the time Google caused significant hand-wringing and soul-searching among institutions traditionally charged with producing or preserving cultural artifacts. (Jeanneney; Venkatraman, 2009) In retrospect, the controversy seems almost quaint in comparison to the current crop of issues – especially the current "disruptions" of established economic models by

DOI: 10.4018/978-1-5225-7659-4.ch038

Uber/Lyft, Facebook, Twitter, Spotify, Snapchat, e-readers, et al. and the encroachments on civil rights via electronic digital surveillance and other intrusions of privacy.

A number of mass-digitization projects have grown in the wake of Google's announcement, including the *HathiTrust*, *Internet Archive*, *Digital Public Library of America (DPLA)*, *California Digital Library*, *Texas Digital Library*, *Gallica*, and *Europeana*. These projects each transcend their roots as localized digital libraries and have simultaneously adapted to and altered the digital landscape. These various MDLs have allowed for and contributed to the ascendancy of our current mass-digitization online culture.

This chapter will describe the characteristics of Massive Digital Libraries (MDLs) and outline their impact upon contemporary information science issues, especially with regard to digital collection metadata, copyright and the diversity of the source collections. Traditionally, libraries have been created to serve particular communities defined by geography, intellectual discipline, or specific end users. However, MDLs in their current trajectories promise–for better *and* for worse—to transcend such limits.

BACKGROUND: DEFINING MASSIVE DIGITAL LIBRARIES

Defining Criteria

Massive Digital Libraries (MDLs) describes a specific class of digital libraries that correspond to the size of a traditional, large brick-and-mortar library. Although other disciplines have discussed digital libraries and archives in terms of computer science, such as in the Very Large Digital Library (VLDL) movement, none have framed the discussion in terms of the principles of library science or the services and content access provided by an actual, working library. (VLDL, 2011)

The following list of characteristics has been proposed to help define MDLs:

- 1. Collection size: surpasses 500,000 texts; prime MDLs comprise tens of millions of texts;
- 2. Acquisitions, collection development & copyright: numerous partnering members contribute content regardless of author or copyright holder permissions and regardless of end-user needs;
- 3. Content type: mass-digitized print books; the resulting searchable digital corpus of texts becomes as important as the individual titles;
- 4. Collection diversity: diversity is dependent upon self-chosen partner members, which can reflect distortions or biases inherent to the source collections;
- 5. Content access: varying degrees of open access exist within MDLs; content is searchable through single, uniform interfaces (search engines & portals) representing all the collections as members of a single entity regardless of source;
- 6. *Metadata: Metadata is gathered and aggregated from multiple sources, with a reliance on new digital description schema;*
- 7. Content / digital preservation: consortium members provide long-term digital preservation strategies at local levels as well as "in the cloud".

These criteria and their attendant issues, though not necessarily unique to digital libraries, require different approaches when dealing with a Massive Digital Library. The issues involved with aggregating millions of previously published print materials into one uniform, yet decentralized, conceptual and online digital space become more complex as size increases. It is important to differentiate MDLs from

11 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/massive-digital-libraries-mdls/215949

Related Content

The Expert's Opinion

Edward J. Szewczak (1998). *Information Resources Management Journal (pp. 35-36)*. www.irma-international.org/article/expert-opinion/51059

Dd

(2013). Dictionary of Information Science and Technology (2nd Edition) (pp. 237-307). www.irma-international.org/chapter/dd/76413

Challenges in Modelling Healthcare Services: A Study Case of Information Architecture Perspectives

George Leal Jamil, Liliane Carvalho Jamil, Augusto Alves Pinho Vieiraand Antônio José Daniel Xavier (2016). Handbook of Research on Information Architecture and Management in Modern Organizations (pp. 1-23).

www.irma-international.org/chapter/challenges-in-modelling-healthcare-services/135759

Virtualization and Its Role in Business

Jerzy A. Kisielnicki (2009). Encyclopedia of Information Science and Technology, Second Edition (pp. 4028-4033).

www.irma-international.org/chapter/virtualization-its-role-business/14180

Infosys Technologies Limited: Unleashing CIMBA

Debabroto Chatterjeeand Rick Watson (2005). *Journal of Cases on Information Technology (pp. 127-142).* www.irma-international.org/article/infosys-technologies-limited/3165