

---

# A VIEW OF THE INTERFACE BETWEEN ETHICS AND METADATA

---

David Haynes (1)

(1) Edinburgh Napier University, Scotland, d.haynes@napier.ac.uk



## Abstract

Despite its advantages in improving access to information, metadata has the potential to cause harm, if used inappropriately. For instance, mass surveillance of phone calls can be used by intelligence agencies to target individuals. Language use can reinforce prejudices; poor language control undermines the efficiency of subject searches; and the opacity of discovery systems reduces the effectiveness of retrieval systems. There are also concerns about who owns the intellectual property associated with metadata creation. The talk concluded with a description of two proposed initiatives: by ISKO to investigate ways of improving metadata use in information discovery systems; and the recruitment for a fully-funded PhD studentship at Edinburgh Napier University to investigate the ethics of metadata.

**Keywords:** Metadata; Ethics; Privacy; State Surveillance; Information Retrieval; Controlled Languages; Information Discovery System.

## 1 Introduction

---

Talk given at the *Transversalidade e Verticalidade na Ciência da Informação* symposium (Transversality and Verticality in Information Science) held at São Paulo State University - UNESP, Marília city, 10-11 August 2023.

It is a great pleasure to be back in Marília for the second time. My first visit was in 2016 and it was a great opportunity to meet colleagues and engage in a dialogue about research, particularly in Knowledge Organization. Thank you for inviting me to join you for this exciting event. I have spent the last few days working with some amazing students as well as with respected research colleagues exploring avenues for joint research and other academic engagements.

The interface between metadata and ethics has been a recurrent theme in my professional and academic life. There were several excellent papers at ISKO 2016 in Rio (Nascimento et al., 2016), ISKO UK in 2019, and others since that have highlighted some of the issues I wish to address today (Haynes, 2018; Haynes & Vernau, 2019). Today's talk will look at the interface between ethics and metadata by considering some specific scenarios. Some ethical responses will be considered before pointing towards current and proposed future research.

I'll start by stating some of the problems associated with metadata.

## 2 Mass Surveillance

---

*"We kill people based on metadata."*

This quote from the former director of the United States of America Central Intelligence Agency (CIA), General Michael Hayden, was widely reported in the press and was the starting point for a panel discussion that I participated in at a Dublin Core Metadata Initiative (DCMI) virtual meeting (Haynes et al., 2021). This quote makes us think about the role of metadata and how it intersects with privacy concerns.

The invasion of privacy by state institutions is not new. Journalists and human rights campaigners have been warning the public for some time. As far back as 1980, Duncan Campbell, an investigative journalist, raised the alarm about mass surveillance by the British intelligence community (Campbell, 2015). In 2013, Amnesty International published a report on the US drone strikes in Pakistan that resulted in scores of civilian deaths (Amnesty International, 2013). Edward Snowden revealed the extent of surveillance of US citizens' communications by the US National Security Agency (NSA) (Greenwald, 2013).

Using metadata is more efficient than monitoring the content of calls. Metadata is structured and provides information on the identities of the callers, who their contacts are, their location, and the timing of their activities. If someone is a person of interest, it is relatively easy to spread the net to capture data about their associates and their contacts. Once an association has been made between an individual and a device, the device's location can be used to target that individual. They might be a terrorist or freedom fighter, depending on your

perspective, or an uninvolved bystander. There is a big difference between arresting a suspect and extra-judicial killing.

### 3 Embedding Prejudice

---

The controlled languages used to describe information resources reflect specific perspectives and worldviews. Holstrom (2022) describes how classification language can be discriminatory and highlights the power imbalance between authorities and user communities.

I started my career as an information scientist abstracting and indexing. I am particularly interested in the way that indexing languages evolve. For instance, when we look at the terminology used to describe indigenous peoples in the Americas (including some of my ancestors), we rapidly get into difficulties:

- AMERICAN INDIANS is used as a term in older literature;
- FIRST NATIONS is widely used as a term in Canada;
- FNMI (First Nations, Metis, Inuit) is also referred to in Canadian literature;
- INDIGENOUS PEOPLE is used as a term around the world, including the United States and Australia;
- NATIVE AMERICANS is still used as a term in the United States of America, and;
- AMERINDS is also used as a term in the Americas.

Words change their meaning depending on the geographical region and evolving views about what is acceptable. Some of these terms overlap and are used interchangeably. There is also a transition between different usages as older terms fall out of favor. Another example is the distinction in Brazil between *índios*– indigenous peoples (related to the Amerinds in the Caribbean) and *indianos*, people from India. Our use of language reflects a particular worldview. Recognizing diversity is good, but labels can also alienate or other people.

### 4 Problems with Searches

---

Metadata permits access to resources based on different characteristics. I'll focus on subject access. One of the principles of a fair society should be citizens' access to information.

That allows people to make informed decisions about their daily lives, what products or services to use, or who they wish to elect to govern their country.

The International Society for Knowledge Organization (ISKO) recently set up a working group to consider some of the issues of lack of access to metadata and its effect on subject retrieval. The premise of that group is that the current generation of discovery systems used in academic and research libraries largely ignore the very rich subject metadata that has been created by specialists. The working group is currently consulting academic and research librarians and it invites participation by members of the knowledge organization community and the wider library and information science community (Haynes et al., 2023).

Controlled vocabularies play an important role in retrieval. An example of the value of controlled vocabularies and applied indexing may help to address well-known problems of spelling variations, ambiguity, and changing use of terminology. In English, there are spelling differences between some words in British and North American usages. If there is no controlled vocabulary, we find that the following examples could give different search results:

- “SULPHUR DIOXIDE” and “SULFUR DIOXIDE”;
- “ALUMINIUM” and “ALUMINUM”.

The other purpose of controlled vocabularies is to disambiguate words so that there is a distinction for instance between the following usages of the word Train:

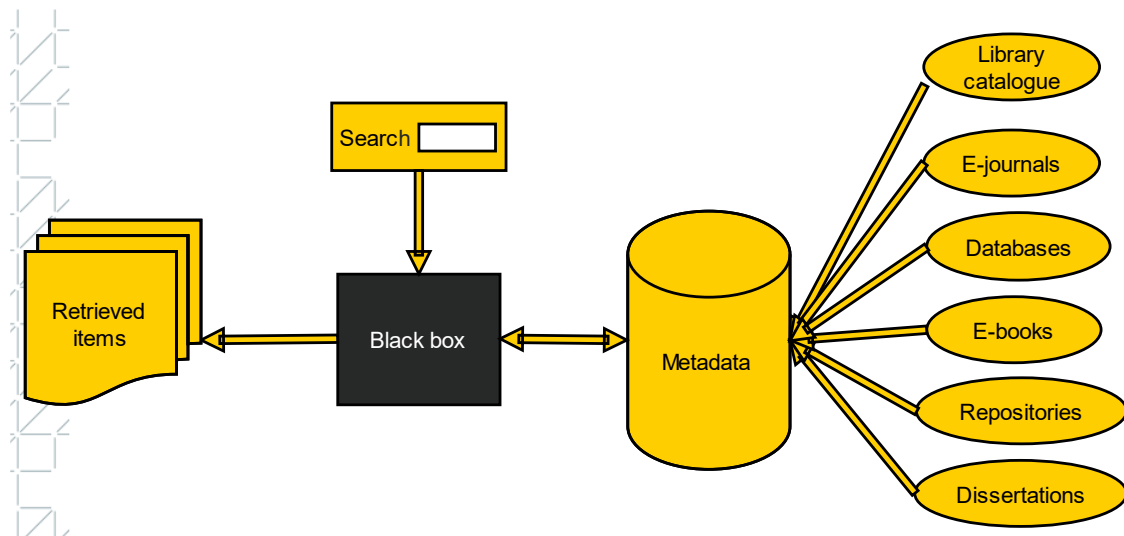
- TRAIN (teaching) – *treinamento*;
- TRAIN (railway vehicle) – *trem*;
- TRAIN (process) – *processo*.

## **5 Opacity of information retrieval systems**

---

Another aspect of metadata usage that can limit access to information is the opacity of search systems. Algorithms for selecting and ranking items tend to be proprietary. They are also deliberately kept hidden to prevent overt manipulation of results. The figure below shows a simplified schematic of a discovery system. We do not know what the criteria are for the selection and ranking of the search results. The lack of transparency makes it difficult to refine and improve searches.

Figure 1 - Opacity of an information discovery system



Source: Author.

This brings us to algorithmic systems. If you create a model based on a limited corpus, you could end up with unintended consequences that reinforce preconceptions and prejudices. Machine Learning has created Large Language Models which can be very helpful tools. However, they are opaque. It is not possible to know how a decision is made by the system. Clearly, many retrieval systems will use metadata, but we are not sure how. This is one of the reasons why European and UK data protection legislation (General Data Protection Regulation, 2016; Data Protection Act, 2018) gives individuals the right to challenge decisions made by AI systems and to require human intervention if requested

## 6 Intellectual property

A final issue of concern is intellectual property. Who owns the metadata? Much of the intellectual effort in creating metadata comes from the authors of papers that are published. There are also professional bodies and national bibliographic authorities who invest a great deal of effort in creating metadata and applying controlled vocabularies. The publishers make this metadata available in various ways. Who owns the intellectual property associated with that metadata, or with a controlled vocabulary? What are the rights of re-use?

## 7 Conclusion

---

All these questions lead us to a consideration of whether metadata is neutral or not. People do not necessarily start with the intent of doing harm, but some well-intentioned actions can lead to harm. Probably most use of metadata does not take ethical considerations into account and this is an area that requires further research. There are many ways in which the library and information profession has responded to some of the ethical challenges that have been outlined. Cutter's (1876) Objects for a Library Catalogue were to: enable finding, show what the library has, and assist in the choice of a book. These still stand and have been substantially incorporated into current IFLA Library Reference Model (Riva et al., 2017) to Find, Identify, Select, Obtain, and Explore. These principles are also reflected in the values and codes of practice of information professionals.

In ISKO we are looking at ways in which retrieval might be improved by encouraging the explicit use of subject metadata in discovery systems. A working group has been set up to describe the problem and create a set of guidelines for the procurement of discovery systems used in research and academic environments. Potential approaches might include a review of interoperability as a way of reconciling different subject vocabularies. Ontology development could also be used as a way of investigating or modeling retrieval systems, to better understand how retrieval decisions are made. Others have used LLMs to design taxonomies (Tang et al., 2023). It should be possible to generate an ontology and thereby expose relationships and inferences. It can help to reveal unexpected relationships that associate crime with specific groups or reflect assumptions about gender roles in the workplace.

ISKO is currently researching the use of metadata in discovery systems and Edinburgh Napier University is offering a fully-funded PhD studentship to investigate the ethics of metadata. There is scope for more work in this area and I hope that this description of some of the ethical issues associated with metadata use will inspire further research.

## References

---

- Amnesty International. (2013). "Will I be next?" US Drone Strikes in Pakistan.
- Campbell, D. (2015). GCHQ and Me. My Life Unmasking British Eavesdroppers. The Intercept. <https://theintercept.com/2015/08/03/life-unmasking-british-eavesdroppers/>.
- Cutter, C. A. (1876). Rules for a Printed Dictionary Catalogue. United States Bureau of Education.
- General Data Protection Regulation, Pub. L. No. EU 2016/679, 78 (2016). <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:02016R0679-20160504&from=EN>.
- Greenwald, G. (2013, June 6). NSA Collecting Phone Records of Millions of Verizon Customers Daily. The Guardian. <http://www.theguardian.com/world/2013/jun/06/nsa-phone-records-verizon-court-order>.
- Haynes, D. (2018). Metadata, Ethics and Trust. Catalogue & Index, 191, 2–4. [https://cdn.ymaws.com/www.cilip.org.uk/resource/collection/C165DDC7-25C3-411B-8137-1BC2A293200B/catalogue\\_and\\_index\\_issue\\_191\\_june\\_2018\\_haynes\\_metadata\\_ethics\\_trust.pdf](https://cdn.ymaws.com/www.cilip.org.uk/resource/collection/C165DDC7-25C3-411B-8137-1BC2A293200B/catalogue_and_index_issue_191_june_2018_haynes_metadata_ethics_trust.pdf).
- Haynes, D., Golub, K., Gnoli, C., Salaba, A., Shiri, A., & Slavic, A. (2023). Improving Search Quality by Enhancing Access to Metadata. Knowledge Organization and Information Discovery: Improving User Experience, Quality and Trust. 7th Biennial ISKO UK Conference. <https://zenodo.org/record/8241636>.
- Haynes, D., Pandit, H. J., & McRae, M. (2021). Metadata and privacy panel. DCMI Virtual 2021. <https://www.dublincore.org/conferences/2021/panels/>.
- Haynes, D., & Vernau, J. (Eds.). (2019). The Human Position in an Artificial World: creativity, ethics and AI in knowledge organization. ISKO UK Sixth Biennial Conference London, 15-16th July 2019. Ergon Verlag GmbH.
- Holstrom, C. (2022). Analyzing the Structure and Dynamics of Control Relationships in the Case of "Illegal Aliens" in the Library of Congress Subject Headings. Knowledge Organization across Disciplines, Domains, Services and Technologies, 133–146.
- Nascimento, F. A., Leite, F. F. J., & Pinho, F. A. (2016). What Gender is this? Challenges to the subject of representation about the gender boundaries. In J. Guimarães, S. O. Milani, & V. Dodebei (Eds.), Knowledge Organization for a Sustainable World: Proceedings of the Fourteenth International ISKO Conference, 27-29 September 2016 (pp. 587–592). Ergon Verlag GmbH.
- Riva, P., Le Boeuf, P., & Žumer, M. (2017). IFLA Library Reference Model. [https://www.ifla.org/files/assets/cataloguing/frbr-lrm/ifla\\_lrm\\_2017-03.pdf](https://www.ifla.org/files/assets/cataloguing/frbr-lrm/ifla_lrm_2017-03.pdf)

Tang, Y., da Costa, A. A. B., Zhang, J., Patrick, I., Khastgir, S., & Jennings, P. (2023). Domain Knowledge Distillation from Large Language Model: An Empirical Study in the Autonomous Driving Domain. <https://doi.org/10.48550/arxiv.2307.11769>.

UK Data Protection Act, (2018).

---

Copyright: © 2023 HAYNES, David. This is an open-access article distributed under the terms of the Creative Commons CC Attribution-ShareAlike (CC BY-SA), which permits use, distribution, and reproduction in any medium, under the identical terms, and provided the original author and source are credited.

---

Received: 17/11/2023

Accepted: 29/11/2023