Digital Policy Hub — Working Paper

# University of Toronto Libraries: A Case Study for AI Governance

## Matthew da Mota

Winter 2024 cohort

## About the Hub

The Digital Policy Hub at CIGI is a collaborative space for emerging scholars and innovative thinkers from the social, natural and applied sciences. It provides opportunities for undergraduate and graduate students and post-doctoral and visiting fellows to share and develop research on the rapid evolution and governance of transformative technologies. The Hub is founded on transdisciplinary approaches that seek to increase understanding of the socio-economic and technological impacts of digitalization and improve the quality and relevance of related research. Core research areas include data, economy and society; artificial intelligence; outer space; digitalization, security and democracy; and the environment and natural resources.

The Digital Policy Hub working papers are the product of research related to the Hub's identified themes prepared by participants during their fellowship.

## Partners

Thank you to Mitacs for its partnership and support of Digital Policy Hub fellows through the Accelerate program. We would also like to acknowledge the many universities, governments and private sector partners for their involvement allowing CIGI to offer this holistic research environment.



## About CIGI

The Centre for International Governance Innovation (CIGI) is an independent, non-partisan think tank whose peer-reviewed research and trusted analysis influence policy makers to innovate. Our global network of multidisciplinary researchers and strategic partnerships provide policy solutions for the digital era with one goal: to improve people's lives everywhere. Headquartered in Waterloo, Canada, CIGI has received support from the Government of Canada, the Government of Ontario and founder Jim Balsillie.

<div style="border: 1px solid green;">

## Key Points

- Existing governance structures for data and information in Canada-based research institutions are varied and often overlapping.

- The gaps in governance that result from this varied system leave vulnerabilities that could be exploited through the implementation of machine learning/artificial intelligence (ML/AI) tools or could leave ambiguities about the right to use certain data or metadata in developing or refining ML/AI models.

- Valuable but underutilized data within research institutions is particularly at risk of being accessed and used by third-party ML/AI tool developers without institutions being properly compensated.

- The sector requires binding standards for ML/AI deployment, alongside broad strategic planning, the promotion of safe experimentation with ML/AI tools and the development of frameworks for institutions to mobilize and exchange their data.

</div>

# Introduction

The beginning of 2024 has seen the release of two ML/AI systems designed for academic libraries[1] as well as one of the first copyright licensing agreements between a licensor and universities specifically outlining restrictions for AI use of copyrighted material.[2] The landscape has shifted from speculation about how AI might be used in the research sector to multiple use cases. Yet the precarity of Canada's governance frameworks for AI has not improved. As the ML/AI tool economy, and the underlying data economy on which it is built, rapidly expand, it is more important than ever to develop effective governance frameworks to protect the research sector in Canada and globally.

The first paper in this series showed that existing governance of AI in research institutions is inadequate, necessitating more comprehensive frameworks. Research institutions support research and are entrusted with preserving human knowledge, making it essential to establish enduring governance. Due to the important societal position of these institutions, improper implementation of AI tools could be catastrophic (da Mota 2024).

Building on the first paper, this second working paper provides a case study for AI governance in a modern research library. The institution, University of Toronto Libraries (UTL), is the largest university library system in Canada[3] and one of the largest in the world. UTL provides an example of the kinds of data and information in research library systems and how they are governed, as well as the risks of ML/AI implementation in the sector.

First, this paper will taxonomize the information and data accessible through the UTL, their custodians, and the laws, contracts and policies that govern access to them. Second, the paper will highlight areas where future conflicts may arise around the implementation of ML/AI systems or where UTL data might be exploited by third parties. Third, the paper will explore existing frameworks that could be adopted in

---

1   See Ex Libris (2024) and the Elsevier Scopus AI search tool at www.elsevier.com/products/scopus/scopus-ai.

2   See Canadian Research Knowledge Network (CRKN) (2024).

3   By volume, Library and Archives Canada is the largest in the country and one of the largest in the world, but it is not exclusively a research library.

a future ML/AI governance model for research institutions in general. Finally, this paper will make recommendations for building resilient ML/AI governance in research institutions.

# Mapping Data and Information Sources in UTL

Like all research libraries, UTL is a confluence of physical collections, digital collections (including born-digital and digitized materials),[4] data and audiovisual collections, special or rare collections and access portals to licensed materials. There is also internal data that is produced through the regular functioning of the library, such as metadata, use data, provenance and donor information and so on. Table 1 shows the groupings of information, where they are held, who can access them and how they are governed.

The different data and information accessible through UTL are governed by various tools: some are more closely governed than others, and most do not currently address the use of ML/AI tools. Governance models for library materials can be categorized into four types: institutional policies, privacy, data governance and copyright. These models exist to varying degrees in the UTL and apply to different materials and data in different ways. The most basic governance structure at UTL is the borrowing policy. The borrowing policy limits access to materials to students, faculty, researchers and members of the public with library cards.[5] Additional terms of use are provided by external platforms to which UTL facilitates access.

## Policies

UTL's Digital Preservation Policy does not directly refer to ML/AI uses but is the most relevant policy on data and information preservation. The policy specifies that "long-term access to digital assets is the purpose of digital preservation," defining digital preservation as a "property of the policies and procedures we use to manage our digital assets," and citing "risk identification and mitigation as a central part of our digital preservation strategy."[6] The focus on long-term access, integrating policy and method and identifying risk, is an effective guide for an ML/AI governance framework in keeping with the Proposed Framework for AI Governance (da Mota 2024, 13).

## Privacy

User privacy in UTL is governed by Canadian law. Privacy governs the personal data of library users, including which materials they access and when, producing identifiable metadata. The relevant pieces of legislation are the Freedom of Information and Protection of Privacy Act (FIPPA)[7] and the Personal Information Protection and

---

4   "Born digital" refers to materials that are produced and published to fit a digital format, such as certain journals or open-access ebooks, in contrast to materials made for print and then later digitized. See Jaillant (2022).

5   See https://onesearch.library.utoronto.ca/loan-services.

6   See https://onesearch.library.utoronto.ca/digital-preservation-and-recovery-services.

7   The basis of FIPPA is the Privacy Act; for more information, see the FIPPA manual (2018) and full legislation: Freedom of Information and Protection of Privacy Act, RSO 1990, c F.31, online: <www.ontario.ca/laws-beta/statute/90f31>.

**Table 1: UTL Information and Data Sources and Governance**

| Category | Owned By | Locations/Types | Access | Governance |
|---|---|---|---|---|
| Physical materials | UTL (stacks and special collections) | Books, periodicals, audiovisual | Catalogue, additional requirements for special collections | Borrowing policies, copyright |
| UTL-owned digital materials | UTL (held on UTL servers) | Books, data sets, periodicals and audiovisual through catalogue | Accessible through catalogue | Borrowing policy, copyright, data governance |
| Licensed materials: public, private and open access | Third parties (Elsevier, Internet Archive, etc.) | Books, data sets, periodicals, audiovisual, etc., in third-party systems | Accessible through catalogue to third parties | Borrowing policy, copyright, licensing, data governance |
| User data (private) | UTL as custodian | Personal and user data | Staff and Google Analytics | FIPPA, privacy policy |
| Internal metadata | UTL | Metadata, linked, bibliographical and bibliometric data | Accessible to staff only through Alma catalogue interface | No governance structure |
| World catalogue | The Online Computer Library Center (OCLC) | WorldCat.org catalogue system | UTL catalogue or worldcat.org | Licensing, terms of use |

*Source:* Author.

*Note*: "OCLC is a global library organization that provides shared technology services, original research, and community programs for its membership and the library community at large. We are librarians, technologists, researchers, pioneers, leaders, and learners" (see www.oclc.org/en/about.html).

FIPPA = Freedom of Information and Protection of Privacy Act

Electronic Documents Act (PIPEDA).[8] As a public institution, UTL is governed by FIPPA,[9] while PIPEDA governs private businesses. Following FIPPA, UTL states it anonymizes user data and does not share identifiable data with third-party organizations.[10] However, they do share anonymized data with Google Analytics that can be identified.[11]

Although PIPEDA does not apply to UTL, third parties such as software providers or content licensors must ensure compliance with PIPEDA where personal data might be exposed. UTL systems collect minimal personal data, storing only what is necessary for basic library functions on University of Toronto (U of T)–owned servers. The main risks to privacy are cyberattacks or the misuse of unwittingly identifiable internet protocol (IP) and domain information shared with Google Analytics. The university makes clear

---

8   See *Personal Information Protection and Electronic Documents Act*, SC 2000, c 5, online: <https://laws-lois.justice.gc.ca/ENG/ACTS/P-8.6/index.html>. Although PIPEDA is not as relevant in this context, the ways it examines user information might play an important role in future discussions of data governance at the level of public institutions, especially in how they engage with third-party service providers.

9   UTL provides a brief outline of FIPPA and its approach to adhering to the law, both in protecting users' personal information and in responding to requests to disclose information: see https://governingcouncil.utoronto.ca/fipp.

10   The UTL statement on online privacy and data collection provides some information about FIPPA and the details of its relationship with Google Analytics, including the explanation that users' data is anonymized but could be identifiable through IP or domain names depending on service providers' naming standards: see https://onesearch.library.utoronto.ca/online-privacy-and-data-collection.

11   Ibid.

that it reserves the right to access information on its systems in some circumstances to ensure safety and consistency of digital services.[12]

## Data Governance

Data governance at UTL varies depending on context. Data in UTL can be broken into available data sets, including research data, proprietary data and metadata. For data generated through research projects, unique data management plans are established for each project.[13] Research involving human data must go through an ethics in human research protocol to ensure ethical and secure governance of human data.[14] Departmental and faculty guidelines for data governance might also shape governance.

Much of the research data at U of T is housed on private or local servers as per established data governance protocols, unless the data is published through one of UTL's open access (OA)[15] repositories such as TSpace[16] or the U of T Dataverse[17] Metadata sources accessible through the UTL catalogue, such as OCLC's World Catalogue, are also governed by a "Services Terms and Conditions" document.[18]

## Copyright

The most widespread and varied governance structure shaping AI/ML use in research institutions is copyright. Every physical and digital publication is subject to copyright laws. Licensed digital materials and their metadata are subject to copyright and to the licensing agreements that govern library access.

The Scholarly Communications and Copyright Office at UTL provides guidance for using copyrighted material in courses, support for publishing and how to use copyrighted material in research.[19] UTL also provides guidance on copyright for generative AI

---

12  The details of the university's right to access information can be found under "Privacy" in the "Appropriate Use of Information and Communication Technology" document from the U of T Division of the Vice-President and Provost: see www.provost.utoronto.ca/planning-policy/information-communication-technology-appropriate-use/.

13  See the Institutional Research Data Management Strategy from the Centre for Research and Innovation Support: https://cris.utoronto.ca/dri_portal/home/.

14  Many of the specific details and guidance documents for data governance through the ethics protocol process are internal documents accessible through the My Research Portal, but some preliminary documentation and advice is on the Ethics in Human Research landing page: https://research.utoronto.ca/ethics-human-research/ethics-human-research.

15  See the U of T Dataverse guidance documents: https://onesearch.library.utoronto.ca/researchdata/getting-started-u-t-dataverse. Borealis is the Canadian repository for Dataverse data (TSpace is the U of T-specific portal). See the terms of use at https://borealisdata.ca/termsofuse/ and preservation plan at https://borealisdata.ca/preservationplan/ from Borealis for more context on the principles of OA data governance.

16  TSpace repository policies and guidelines for how data and research are published and on what terms can be found at: https://tspace.library.utoronto.ca/about/collectionpolicies.jsp.

17  To understand more about data governance in OA contexts, see the Findability, Accessibility, Interoperability and Reusability (FAIR) framework for data governance. The FAIR principles are foundational thresholds for data governance that can form the basis of any governance model assisting in "the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals" (Wilkinson et al. 2016).

18  See www.oclc.org/content/dam/ext-ref/worldcat-org/terms.html.

19  See the UTL Copyright Office website at: https://onesearch.library.utoronto.ca/copyright/scholarly-communications-and-copyright-office; for guidance on the use of licensed content, see https://onesearch.library.utoronto.ca/copyright/using-library-licensed-content.

(GenAI), identifying the uncertainty of copyright status for its training materials and that the material produced with GenAI is not copyrightable as potential risk vectors.[20]

Every licensor of content accessible through the UTL has a licensing agreement based on the CRKN's Model License or Open Access Model License (CRKN 2016, 2023a).[21] The model licences provide basic requirements agreements and give insight into the agreed-upon thresholds for fairness and access to information by Canadian universities. CRKN's licensing principles are sustainable scholarly communications, equity of access, OA scholarship and transparency.[22] CRKN advocates for open research and open science principles in Canada, high standards for metadata and persistent identifiers adhering to the FAIR[23] data governance principles and respecting Indigenous data through engagement with CARE[24] and OCAP®[25] principles[26] (CRKN and Canadian Association of Research Libraries 2023).

## Metadata and Licensing

Metadata is slated to become a central concern of libraries seeking to adopt ML/AI tools as metadata is essential for refining and developing ML/AI models. However, there has been little governance of this issue thus far. The current CRKN model licence (as of March 2024) mentions data mining but not ML/AI uses: "Members and Authorized Users may apply automated tools and processes to the Licensed Materials for the purposes of textual analysis and visual mapping of textual and/or statistical relationships within the context of scholarship, research, and other educational purposes."[27] This clause would likely be interpreted to exclude its use in any for-profit or competing ML/AI tools. There are several renegotiations of licensing agreements due in 2023–2024 that may add ML/AI language to the model licence.[28] This year, CRKN negotiated an amendment to their agreement with Elsevier on behalf of their member organizations,[29] which restricts use of licensed materials by institutions using AI/ML tools. The section concerning AI use states that subscribed content cannot be used to "create a competing commercial product," to "adversely disrupt the functionality of the Subscribed Products" or to "reproduce or redistribute the Subscribed Products to third party artificial intelligence

---

20   For UTL guidance on copyright and AI use, see https://onesearch.library.utoronto.ca/copyright/generative-ai-tools-and-copyright-considerations. There is also a useful flowchart from the University of Alberta that outlines copyright status for all types of material: www.ualberta.ca/faculty-and-staff/copyright/intro-to-copyright-law/licensed-royalty-free-content/pd-flowchart---types.html.

21   The CRKN represents most Canadian universities in advocacy for OA principles in education and research, policy development and negotiations for licensed content. The CRKN's model licence agreements are important tools for unifying university negotiations with large licensors: see CRKN (2016, 2023a).

22   See www.crkn-rcdr.ca/en/crkn-licensing-principles.

23   See note 16.

24   CARE is a set of principles for Indigenous data governance, which stands for Collective Benefit, Authority to Control, Responsibility and Ethics: see www.gida-global.org/care.

25   OCAP® is a registered trademark of the First Nations Information Governance Centre. More details on the OCAP® framework can be found at: https://fnigc.ca/ocap-training/.

26   See CRKN (2023b).

27   See www.crkn-rcdr.ca/sites/default/files/2024-04/CRKN-Elsevier_2024-2026_License%20Amendment_2024-03-28_red.pdf .

28   For more on current licensing negotiations by CRKN, see www.crkn-rcdr.ca/en/license-negotiations.

29   Elsevier is one of the largest copyright owners and publishers in academic publishing sectors worldwide, and they are one of two companies with newly available AI/ML-enabled tools for use in academic libraries.

tools" with exceptions for projects that are for research only.[30] This agreement will likely become a model for future agreements to protect licensed materials, as CRKN could add this language to their model licence. The amendment sets a favourable precedent for universities since there is no language requiring universities to allow the licensor to access institutional data.

## Ungoverned Data

The final grouping of data in the UTL, internal proprietary, lacks a governance structure, and this data is often unstructured and unlabelled. It is valuable to institutions and third parties because they are essential for creating ML/AI tools, necessitating clearer governance. Having effective data governance structures not only protects data but also provides clear paths to mobilizing, labelling and using unstructured or unused data. There is a significant opportunity cost to not mobilizing this data, as it could provide insight into the functions of the library system and connections between materials. Many of these data sets are unstructured because they pertain to physical library collections in which only some of the metadata has been digitized. The ungoverned data ranges from unstructured data sets produced through library use to data contained within physical collections such as optical character recognition data from undigitized texts. For example, the UTL's Thomas Fisher Rare Book Library recently held a special exhibition of data visualizations represented in their physical collections, containing tables, graphs and data that have never been digitized.[31]

# Analysis

The existing governance structures in UTL cover basic privacy, data protection, copyright and security concerns that are relevant to the ML/AI context. However, there are several gaps and vulnerabilities that arise with ML/AI tool deployment. The primary vulnerabilities come from limits in norms for anonymization and privacy, ungoverned data sources and the way that copyright licensing and control of content opens the sector to anti-competitive practices and security risks.[32]

Regarding privacy, UTL does well to limit and anonymize user data. If a future ML/AI system was deployed for analytics, this could require third-party access to user data (such as the current Google Analytics arrangement) and would need attention to ensure FIPPA and PIPEDA are adhered to by all parties. UTL could also take steps or provide guidance to ensure that Google Analytics and other third parties cannot identify users through IP or domain names. A less immediate but more significant risk comes from speculation that ML/AI tools may render current anonymization ineffective through identifiers not anticipated by current privacy protocols,[33] thereby requiring a higher

---

30   See www.crkn-rcdr.ca/sites/default/files/2024-04/CRKN-Elsevier_2024-2026_License%20Amendment_2024-03-28_
     red.pdf, Appendix 1, for the full text of the AI clause from the CRKN-Elsevier licence amendment.

31   See https://fisher.library.utoronto.ca/exhibition/emerging-patterns-data-visualization-throughout-history.

32   This is a broader issue in libraries and academic publishing in general, which has been identified and discussed for
     decades, and organizations such as CRKN have been advocating for pro-research policies and agreements to resist anti-
     competitive and monopolistic practices. See notes 19, 20, 24 and 25 for more on CRKN's work in this area.

33   An article from Imperial College London summarizes the key risks of ML/AI use to deanonymize: see Brogan (2019). Two
     further articles explore how anonymization might be adapted to this risk: see Thompson (2023) and Rustad (2024). A
     blog post on AIhub identifies the risks of irreversible anonymization: see Moreton and Jaramillo (2021).

standard of "privacy by design."[34] UTL and U of T have robust standards for data governance in research projects, but these standards are not applied to other data, namely, metadata and unstructured data in the UTL system. Licence agreements and terms of use for data libraries provide some guidance on how data and metadata can be used, but internal UTL data remains ungoverned. Does UTL providing metadata to the World Catalogue allow OCLC to use that metadata in developing ML/AI tools? Does the fact that UTL publicly displays the metadata for licensed material on its catalogue mean that UTL can use that metadata to train their own ML/AI tools? There currently are no clear answers to these questions.

The copyright landscape does not currently provide guidance for how ML/AI can use materials. In recent cases, it seems that training materials of models such as ChatGPT violated copyright.[35] Although some claim using copyrighted material for training is fair use,[36] if it is later determined that this violates copyright, it would force developers to pay fines and licensing fees or retrain their large language models (LLMs), causing significant disruption and financial loss. Without extensive precedent for this issue, the outcome is uncertain and may not increase security or clarity. Regarding essential metadata for developing ML/AI-powered library tools, it is unclear whether copyright applies. However, if licensors restrict metadata through agreements, those contracts may functionally supersede copyright as they will define how licensed materials are used by institutions.

Unstructured data in UTL and other library systems do not have clear governance models. With these valuable assets ungoverned, the potential for exploitation or misuse must be addressed. One major risk is companies that license materials to libraries negotiating licensing agreements that allow them access to all data within a system such as UTL. This would be something a university could consent to in negotiations that would erase the data's value as an asset. Universities may rush to have the latest capabilities and preserve continued access to the latest journals, leading them to sign contracts that would undermine future opportunities. As the internal data in research institutions is an asset, there needs to be clear frameworks reflecting the data's ownership by the institution, including methods for transacting or sharing the data to the institution's benefit with integrated data preservation methods.

## Anti-competitive Practices and Leveraging Data Value

A recent discussion paper from the Competition Bureau of Canada explores the potential for anti-competitive and monopolistic environments in the AI economy. The report identifies data supply as one of the three main markets for AI (Competition Bureau of Canada 2024, 8), suggesting that "publicly available data...could be fully exhausted in the next few years" (ibid., 11).[37] The report also speculates on whether proprietary data

---

34  The International Association of Privacy Professionals provides a resource list for various perspectives on privacy-by-design methods and frameworks: see https://iapp.org/resources/article/oipc-privacy-by-design-resources/.

35  OpenAI is quoted as saying that it is "impossible" to train LLMs without copyrighted materials. See Edwards (2024) and Milmo (2024) for context.

36  The fair use argument is central to the defensive claims of OpenAI and other AI companies regarding copyright violation, detailed here in analysis from the Association of Research Libraries: see Klosek and Blumenthal (2024).

37  The report cites Pablo Villalobos et al.'s article "Will we run out of data? Limits of LLM scaling based on human-generated data," which suggests that high-quality public data will be exhausted as early as 2026, while low-quality public data (language and image data) will endure much longer (Villalobos et al. 2022).

will have a market (ibid., 12); however, in the research sector, proprietary or internal data will play a significant role in building and training ML/AI research tools. The report suggests that access to data and computers may be barriers to entry for new or small firms, especially if that data is proprietary (ibid., 14). Finally, the report identifies risk of anti-competitive practices through predation (incurring losses to gain dominance in the market), exclusion (preventing competition by preferencing their own products) and tying or bundling of services (ibid., 17–18). The few large firms that supply both library services and licensed materials could bundle ML/AI services, library services and licensed materials in order to limit competition.

The varied landscape of data governance in UTL and other university library systems makes these systems vulnerable to exploitation by large library services and copyright licensing companies. Companies such as Clarivate and their subsidiaries own vast amounts of content while also offering library management software and now ML/AI-powered services for libraries.[38] This creates an anti-competitive and monopolistic environment where a few US-based companies, own most of the scholarly material and the systems needed to manage that material.[39] As ML/AI tools are driven by data, this monopoly may bar smaller companies or institutions from creating ML/AI tools due to a lack of data, especially if proprietary institutional data is also controlled by larger firms through contracts or other means. Leveraging data value will be the point on which the entire ML/AI economy turns. This is particularly significant in the research sector where industry has the power to dictate terms to institutions. Leveraging high-quality data in research institutions is key to balancing that dynamic for the independence of research institutions.

An anti-competitive ML/AI services market for research institutions also poses a significant risk to research security. In its most publicized form, research security in Canada has focused on military research internet protocol.[40] However, research security extends to influence or interference with research, which can have long-term effects on research freedom. The fact that all major owners of library software services and academic copyright are in the United States poses a research security risk, despite the United States being Canada's ally, and opens the sector to vulnerabilities from other external forces that might exploit weak copyright and data laws. Having power concentrated in only a few firms promotes an uncompetitive environment and threatens research independence. This is an important consideration when thinking of the long-term health of any nation's research and development landscape.

---

38  Ex Libris Group is owned by Clarivate. See the description of Ex Libris's AI Enrichment Metadata Generator: https://exlibrisgroup.com/announcement/ai-enrichment-metadata-generator/.

39  Clarivate owns ProQuest (one of the largest copyright owners in the world), which in turn owns Ex Libris, the largest provider of library management software services; among this software is Alma, which UTL uses for its systems.

40  The Government of Canada provides guidance documents and training for researchers and institutions to improve research security: www.canada.ca/en/services/defence/researchsecurity.html. Several large Canadian research institutions, including U of T, have opened research security offices under the Division of the Vice-President, Research and Innovation: https://research.utoronto.ca/safeguarding-research/safeguarding-research.

# Implementation and Development

The first publicized GenAI tool for cataloguing is Ex Libris's AI Enrichment Metadata Generator,[41] which generates metadata and catalogues texts. How the governance of libraries' data would function in that kind of system is unclear and requires further scrutiny to assess whether it would be a secure solution. If the ML/AI models from Ex Libris and other providers exist inside individual library systems, with libraries maintaining their own data,[42] then this tool provides the benefits of large language models while limiting the risks. However, if the metadata generator requires proprietary data from libraries, this could erode proprietary data value.

The only AI-enabled library search tool currently being marketed is Elsevier's Scopus AI, which is "an intuitive and intelligent search tool powered by generative AI" that provides enhanced search results from the Elsevier catalogue for materials from 2013 and on.[43] Scopus AI claims to cite referenced texts, provide article summaries, suggest "go deeper" questions, generate concept maps and more.[44] Scopus AI could yield positive results for research institutions, and it currently does not appear to require institutional data. Given Elsevier's mostly self-protective, non-predatory licence amendment with CRKN, it seems they are currently focusing on AI/ML tools that manage their own catalogues rather than institution-wide catalogue search tools, thus posing less risk to institutional data.

## Potential Governance Models

In this emerging space, any ML/AI governance models must begin with strong data governance that includes technology governance standards[45] and integrated mechanisms in the data sets to ensure consistency and compliance. There are some existing tools and models that would serve as the foundation for a UTL data governance model for ML/AI, which could work across the Canadian research sector.

### Core Principles

The UTL Digital Preservation Policy[46] provides a baseline for preservation, long-term access and risk identification. These are important and proactive strategic principles for a strong governance model that could be combined with other principles such as the Proposed Framework for AI Governance (da Mota 2024, 13).

---

41  See https://exlibrisgroup.com/announcement/ai-enrichment-metadata-generator/.

42  For example, Canadian company Cohere AI (see https://cohere.com/) provides GenAI services that can be adapted to specific use cases within businesses. Cohere does not have a specific tool for research institutions but could be in the space if copyright and data control issues could be addressed in the library services sector.

43  See www.elsevier.com/products/scopus/scopus-ai.

44  Ibid.

45  Technology governance standards are guidance documents, drafted by stakeholders in a certain area, that govern technical and practical elements of technology use in a certain sector or context. Standards are developed by certified bodies that facilitate stakeholder engagement and standard drafting and execute regular assessments to ensure the standard is continuing to meet its goals and that all adherents to the standard are following it.

46  See note 6.

Another important consideration for a governance model will be classification systems and ontologies for data labelling.[47] Having interoperable systems for data labelling that integrate the main principles of data security and preservation will ensure data governance longevity. The United Nations Educational, Scientific and Cultural Organization (UNESCO) and other organizations have extensive resources for key terms, labelling and lexicons in various disciplines, which could support interoperability.[48] The use of persistent identifiers (PIDs) for digital content would also help this effort by making digital materials, and individual researchers, more findable through a stable classification system.[49]

## Data Maturity Models

If there is a standard set of guiding principles and a classification system, then the next step for a successful data governance model at UTL and between institutions would be to integrate a data maturity model. Data maturity models provide a standard against which institutions can measure their data sophistication for individual data sets and across the institution. There are many examples of data maturity models,[50] but Mark S. Fox, Bart Gajderowicz and Dishu Lyu provide a useful framework for an academic context. Their model has six levels: level 1 comprises descriptive, temporal and geospatial identification (Fox, Gajderowicz and Lyu 2024, 23). Level 2 "focuses primarily on access" (ibid.). Level 3 includes "additional content and versioning information" and incorporates FAIR principles (ibid., 24). Level 4 identifies if there are individual data or Indigenous data in the set (ibid., 25).[51] Level 5 "focuses solely on capturing additional meta-data related to FAIR principles" (ibid., 27). Finally, level 6 "provides basic statistics and data quality" (ibid., 20). This model can be implemented through a plug-in in a library or data management system with user-friendly inputs (ibid., 28).

## Promoting Innovation

While there are existing models and frameworks that could be adapted into an ML/AI governance model for higher education and research, policies and technical frameworks can only go so far if the institutions do not have the capacity or culture to experiment with and implement these technologies. UTL and U of T have some experimentation in ML and AI currently through several initiatives.[52] However, scholars and staff who could make use of tools to analyze proprietary data are not always equipped to experiment

---

47  A data ontology is "a framework to represent information, and as such it can be representationally successful whether or not the formal theory used in fact truly describes a domain of entities" (Vázquez 2018).

48  See the UNESCO thesaurus: https://vocabularies.unesco.org/browser/thesaurus/en/.

49  U of T uses PIDs for digital materials and the ORCID ID platform which provides PIDs for individual researchers to make their work more shareable: see https://onesearch.library.utoronto.ca/researchdata/orcid. Find more details on PIDs for texts in libraries at: https://libguides.lib.umanitoba.ca/RI-and-Profiles/pids.

50  See the Microsoft Responsible AI Maturity Model (Vorvoreanu et al. 2023).

51  The FAIR data management principles and the OCAP® principles for Indigenous data are both mentioned in the data maturity model and could provide guidance for a comprehensive data governance model supporting broader ML/AI governance for research institutions.

52  U of T's BMO Lab explores uses of AI in creative practice and performance explicitly separated from for-profit AI research that may "go to market," providing an interesting space for pure experimentation of the transdisciplinary implications of the research: see https://bmolab.artsci.utoronto.ca/.  Professor Paolo Granata proposed a course taught by AI tools, integrating the technology with classic pedagogy to experiment with the implications of AI use in education (Faculty of Arts & Science 2023). However, the project was cancelled due to uncertainty about its implications.

with ML/AI tools. Connecting technical specialists with research specialists seeking to solve problems with ML/AI would promote innovation.

By mobilizing existing funds, programs and resources to promote ML/AI experimentation through low-risk exploratory projects, the institution could develop a larger, transdisciplinary community of practice. Ryan Cordell (2022) argues that while librarians have been discussing the potential impacts of AI for years, there is a hesitancy to experiment. He sees low-risk, low-scale experimentation as an essential step to learning how AI can be used in libraries and research (ibid.). Such projects in any university library could benefit from additional funds, but it is not a requirement as institutional and governmental efforts already exist that support small projects.[53] The main push that such an effort would need is administrative buy-in at an institutional level, a central document or space advertising existing funds and initiatives to apply for, and institutional policy and strategic thinking to encourage innovation. These minor efforts would build a culture of "exothermic innovation" (Paprica 2023),[54] which would be self-sustaining and would show the best paths for ML/AI use and deployment in each institution. A robust culture of experimentation could help test what works and what does not, as well as build norms and practices that inform policy development. This is also an essential effort for mobilizing the ungoverned data through selectively developing tools for specific use cases.

## Public Data Infrastructure

Another essential piece of the puzzle of ML/AI deployment in the research sector is establishing mechanisms and markets for institutions to monetize, share and exchange data sources between institutions. If managed well, the value in institutions' internal proprietary data could be immense. Here, the word "value" refers both to monetary value and the potential value in supporting research, especially when considering opportunities for pooling or connecting data sets across institutions to amplify their scope.

There are several interesting models for data exchange platforms that might provide insight into what a Canadian or international research data exchange could look like. First, systems such as TSpace, Borealis and Dataverse provide models for effective management, preservation and access to large repositories of data at various scales. Building a data exchange infrastructure into the existing OA data platforms such as Borealis might be a useful way to leverage existing tools in order to maximize the

---

53  The LEAF+ Generative AI in Teaching and Learning project offers funding for AI research in teaching and learning (https://ocw.utoronto.ca/leaf-ai/). The Social Sciences and Humanities Research Council's institutional grants provide funding to small-scale initiatives in post-secondary institutions, along with their suite of larger research funds: see www.sshrc-crsh.gc.ca/funding-financement/programs-programmes/institutional_grants-subventions_institutionnelles-eng. aspx. The Instructional Technology Innovation Fund (https://usc.utoronto.ca/service/instructional-technology-innovation-fund-itif/) provides funds for specific technology implementation or projects in the university, alongside other UTL and U of T project-funding opportunities. Finally, the Engineering Strategies and Practice at U of T gives first-year engineering students the chance to work on a real project to solve a problem in the community or university, which could be used to connect engineering students with librarians and researchers looking to experiment with ML/AI solutions (www.engineering.utoronto.ca/engineering-strategies-practice/).

54  Alison Paprica uses an analogy from chemistry to argue that innovation can be exothermic (needing low amounts of energy going in to sustain higher energy coming out, like lighting wood on fire) and self-sustaining, or endothermic, and therefore resource intensive (needing high levels of energy to sustain a reaction, like nuclear fusion). Paprica (2023) frames this as a model for exploring the best paths for innovation in which following areas of exploration that generate more ideas and take fewer resources and less energy to sustain can be a way to identify sustainable innovation, while also recognizing that sometimes high-cost innovations with less of a cascading effect might be necessary.

efficiency and reach of a data exchange platform for research. A platform could be focused on exchanges between institutions and could be developed into a marketplace to connect research institutions with industry, government or other sectors seeking specialized data sets. It may also be important for a platform to allow institutions to hold their data on their own servers and facilitate access through Application Programming Interface rather than selling the data sets as a package to be downloaded.[55]

When thinking of a national data market, the Chinese system of data exchanges (regional, municipal and national levels) is the only current example of a large network of spaces and platforms in which data can be bought, sold and traded.[56] The Chinese model could provide insight into large-scale data mobilization, valuation of data and platform configurations.

Moving from the national scale to non-governmental examples, the now defunct Innovate Cities Data Trust by CityShield was a not-for-profit platform that sought to provide a venue for public, private and academic partners to share data for mutual benefit while ensuring a high standard of data protection, initially in the smart cities space but then subsequently moving into other areas.[57] CityShield's Chief privacy officer was Ann Cavoukian (2009), who is the information and privacy commissioner of Ontario and the creator of the concept of privacy by design.

Gaia X is another platform collaborating across public, private and academic spaces that provides data exchange services and data governance while maintaining the sovereignty of the data owner, maximizing the benefits of data sharing while limiting the risks and eliminating the need to sell off or give up rights to owned data.[58]

A recent article in *The Globe and Mail* presented a very different argument that Canada should develop a sovereign wealth fund for data to pool Canadians' private data in a powerful trust that could yield great benefits for the nation, while also protecting data from being siphoned for free by large multinationals (Birch 2024). The need to disclose personal data for this proposed fund might not be appealing to Canadians, but the idea of pooling resources to leverage value at a national level is useful.

The main thread connecting these examples is the need for nations and sectors to have better government- or sector-run data infrastructure. Efforts by institutions to leverage their data value through individual contract agreements with service providers would likely prove futile. However, if institutions work together in the interest of research and education, ensuring fair compensation for assets, they could have a chance at determining their own ML/AI implementation path and not rely so heavily on foreign monopolies.

---

55 The API access model, rather than the buy-and-download model for data access, is increasingly the norm for data exchange as companies seek to hold on to their own data and simply supply access to users for security and data sovereignty reasons. Developing APIs for data access will be an essential part of digital public infrastructure (DPI) development over the coming years. For more on DPI standards development, see G20 (2023).

56 See Alex He and Robert Fay (2022) on digital governance in China, He (2023) on China's state-centric data governance model, and He and Rebecca Arcesati (2023) from the IARIW-CIGI conference on the value of data, for more on China's national data-trading system.

57 See https://cityshield.ca.

58 See https://gaia-x.eu/.

## Standards and Non-technical Governance

The research sector is a complex web of interdependent institutions and interests from private, public and not-for-profit sectors. New governance structures for AI in research must be adaptive to this diverse landscape and must engage stakeholders across the spectrum of interests represented in the area. Governance standards are one important piece of the governance puzzle (along with internal policies and national legislation) because of their ability to effectively represent and engage multiple interests at once. Recent criticism of emerging standards regimes for AI correctly identifies that standards are not the apolitical, objective documents some might claim as they are constructed by stakeholders with political interests, potentially undermining "the functional utility of standards as soft law regulatory instruments" (Solow-Niederman 2024). However, despite the potential for bias, standards are an effective and entrenched tool for soft law regulation, which should be used in confluence with legislation and policy, while supporting efforts to make the standards development process as democratic and broadly engaging as possible.

The Canadian Digital Governance Standards Institute (DGSI) currently has a new standard proposal for ML/AI implementation in research institutions under consideration, which will seek to establish benchmarks for reliable and safe ML/AI implementation.[59] Having a specific standard for universities and research institutions will help to establish level ground before institutions begin adopting and developing ML/AI tools more extensively.

A relevant international standard that might support the DGSI's efforts is the International Organization for Standardization's general AI governance standard, ISO 42001,[60] which established basic governance structures for AI implementation in enterprises in general. ISO 42001 does not address many of the very specific concerns in the research sphere, but it does provide a solid base on which to build a standard specific to universities and research institutions.

Regarding some of the uncertainty about privacy in the ML/AI context, the Canadian Information Privacy Protection Framework standard (CAN/DGSI 109-2)[61] may provide additional steps to ensure the protection of private data within research institutions. Individual universities could sign on to these standards individually, but it might also be effective for government to make certain standards a requirement as a condition to ensure ML/AI safety.

The Standards Council of Canada (SCC) published their Data Governance Standardization Roadmap, which outlines the various standards-based efforts necessary to build a resilient data landscape for Canada's future (SCC 2021). The road map makes 35 recommendations for standardization, including establishing accountability frameworks, certification for emerging roles, digital literacy, cybersecurity, privacy and data management governance (ibid., 24). This road map should be essential reading

---

59  The standard proposal can be found on the DGSI website under the "Exploratory Work" section of the "Find a Standard" page: https://dgc-cgn.org/standards/find-a-standard/.

60  See www.iso.org/standard/81230.html.

61  See https://scc-ccn.ca/standards/notices-of-intent/cio-strategy-council/canadian-information-privacy-protection-framework.

within research institutions to help flag some of the standardization efforts that can be undertaken at an institutional level, as well as nationally and internationally.

The concept of essential reading for data governance raises the important topic of having accountable and responsible stewards in place who can help to develop and implement policies within and between institutions. Some universities have created chief AI officers[62] or steering committees on AI, but these roles must be oriented in such a way that they can identify key issues within their research library systems, establish governance frameworks and ensure these frameworks are implemented consistently. Whether this role requires a new position or is made part of an existing position will depend on individual institutions. However, the key point here is that there must be a person or group of people responsible for policy to be effective in consort with national standards development and legislation.

The need for data infrastructure, data exchange capacity, digital standards development, experimentation and innovation, protection of research independence and maintaining a high quality of information and data preservation all support a coherent strategy for ML/AI implementation in research institutions. While many of these efforts can be explored on an institutional or inter-institutional level, there needs to be a broad vision and coherent strategy from planning to implementation to ensure high-quality research and information systems in Canada and around the world. It is necessary to have a clear strategy with tangible tools and methods for implementation led by government, along with an open discussion among citizens about what the future of AI should be in Canada and how it should (and should not) impact essential systems and institutions, particularly regarding the research on which all other sectors and areas of the economy have built their foundations.

# Recommendations

- Develop a guidance standard for ML/AI implementation in research institutions. This should consider, and integrate where possible, elements of FAIR, CARE, OCAP®, data maturity, OA principles, the Proposed Framework (da Mota 2024) and other relevant frameworks in consultation with all stakeholders.

- Promote low-risk, solution-oriented experimentation with ML/AI in libraries to build norms, develop tools and reduce reliance on external companies.

- Mobilize and monetize data for the benefit of institutions:

  - Creation of a data exchange platform for research institutions to exchange and unlock the value of their data.

  - Pooling data and resources across institutions to enhance the impact on research and leverage data value in negotiation with third parties.

---

62 For more on Western University's chief AI officer, see Ferguson (2023).

- Create a federal/provincial government-led working group to lead a discussion on strategy and long-term planning for ML/AI use and the future of research in Canada.[63]

## Acknowledgements

I would like to thank my supervisor, Paolo Granata, as well as Paul Samson, Aaron Shull, Mai Mavinkurve and my peer-reviewers, Ozan Ayata and Mahatab Uddin, for reading and providing notes on this working paper. I would also like to thank Reanne Cayenne and Dianna English for their extensive support of my work and of the Digital Policy Hub.

## About the Author

Matthew da Mota is a post-doctoral fellow at CIGI's Digital Policy Hub, where he researches the uses and governance of artificial intelligence and large language models within universities and public research institutions (including research libraries and archives). He will also explore the implications of these public institutions' research policies in the private sector and in government. Matthew's other research interests include the connections between history and imperialism, the creation and uses of historical narratives over time, the philosophy of history, propaganda, and the way that media and technologies shape the circulation and preservation of information. He recently obtained his Ph.D. from the University of Toronto's comparative literature program.

# Works Cited

Birch, Kean. 2024. "Canada needs a sovereign wealth fund — built by monetizing our personal data." *The Globe and Mail*, March 21. www.theglobeandmail.com/business/commentary/article-canada-needs-a-sovereign-wealth-fund-built-by-monetizing-our-personal/.

Brogan, Caroline. 2019. "Anonymising personal data 'not enough to protect privacy', shows new study." Imperial College London, July 23. www.imperial.ac.uk/news/192112/anonymising-personal-data-enough-protect-privacy/.

Cavoukian, Ann. 2009. "Privacy by Design: The 7 Foundational Principles." Toronto, ON: Information and Privacy Commissioner of Ontario. www.sfu.ca/~palys/Cavoukian-2011-PrivacyByDesign-7FoundationalPrinciples.pdf.

Competition Bureau of Canada. 2024. "Artificial intelligence and competition: discussion paper." March. Gatineau, QC: Competition Bureau Canada. https://publications.gc.ca/site/eng/9.935280/publication.html.

Cordell, Ryan. 2022. "Closing the Loop: Bridging Machine Learning (ML) Research and Library Systems." *Library Trends* 71 (1): 132–43. https://doi.org/10.1353/lib.2023.0008.

---

63 The Natural Sciences and Engineering Research Council, the Social Sciences and Humanities Research Council, the National Research Council, the Office of the Chief Science Advisor and Innovation, Science and Economic Development Canada would all be important collaborators for this kind of office. The proposed Canadian AI Safety Institute, announced by Deputy Prime Minister Chrystia Freeland (2024), might also be a useful collaborator or leader in this effort.

CRKN. 2022. "Model License." Ottawa, ON: CRKN. www.crkn-rcdr.ca/sites/default/files/2022-07/CRKN%20Model%20License_2022_FINAL_EN.pdf.

— — —. 2023a. "Open Access Model License." Ottawa, ON: CRKN. www.crkn-rcdr.ca/sites/default/files/2023-04/CRKN%20OA%20Model%20License_2023_FINAL_EN.pdf.

— — —. 2024. "CRKN Signs Read-and-Publish Agreement with Elsevier." Press release, April 8. www.crkn-rcdr.ca/en/crkn-signs-read-and-publish-agreement-elsevier.

CRKN and Canadian Association of Research Libraries. 2023. "Towards Open Scholarship: A Canadian Research and Academic Library Action Plan to 2025." May. www.crkn-rcdr.ca/sites/crkn/files/2023-05/Towards%20an%20Open%20Scholarship%20Action%20Plan_EN_FINAL.docx_0.pdf.

da Mota, Matthew. 2024. "Toward an AI Policy Framework for Research Institutions." Digital Policy Hub Working Paper. www.cigionline.org/publications/toward-an-ai-policy-framework-for-research-institutions/.

Department of Finance Canada. 2024. "Remarks by the Deputy Prime Minister on securing Canada's AI advantage." Government of Canada, April 7. www.canada.ca/en/department-finance/news/2024/04/remarks-by-the-deputy-prime-minister-on-securing-canadas-ai-advantage.html.

Edwards, Benj. 2024. "OpenAI says it's 'impossible' to create useful AI models without copyrighted material." *Ars Technica*, January 9. https://arstechnica.com/information-technology/2024/01/openai-says-its-impossible-to-create-useful-ai-models-without-copyrighted-material/.

Ex Libris. 2024. "Introducing our AI Generated Metadata." *Artificial Intelligence Blog Series*, February 25. https://exlibrisgroup.com/blog/artificial-intelligence-blog-series-introducing-our-ai-metadata-generator/.

Faculty of Arts & Science Staff. 2023. "U of T prof to offer experimental course taught with AI tools like ChatGPT." *U of T News*, April 10. www.utoronto.ca/news/u-t-prof-offer-experimental-course-taught-ai-tools-chatgpt.

Ferguson, Keri. 2023. "Western appoints Mark Daley as first-ever chief AI officer." *Western News*, September 27. https://news.westernu.ca/2023/09/western-appoints-daley-chief-ai-officer/.

Fox, Mark S., Bart Gajderowicz and Dishu Lyu. 2024. "A Maturity Model for Urban Dataset Meta-data." *arXiv*, February 23. https://doi.org/10.48550/arxiv.2402.05211.

G20. 2023. "Annexure 1: G20 Framework for Systems of Digital Public Infrastructure." Digital Economy Ministers Meeting: Online Document & Chair's Summary, Bengaluru, India, August 19. https://g7g20-documents.org/fileadmin/G7G20_documents/2023/G20/India/Sherpa-Track/Digital%20Economy%20Ministers/2%20Ministers%27%20Annex/G20_Digital%20Economy%20Ministers%20Meeting_Annex1_19082023.pdf.

He, Alex. 2023. *State-Centric Data Governance in China*. CIGI Paper No. 282. Waterloo, ON: CIGI. www.cigionline.org/publications/state-centric-data-governance-in-china/.

He, Alex and Rebecca Arcesati. 2023. "Better Governance to Unleash the Value of Data: China's Practice of Building a Data Trading System." Paper presented at the IARIW-CIGI Conference on the Valuation of Data, Waterloo, ON, November 2–3. https://iariw.org/wp-content/uploads/2023/10/He-Arcesati.pdf.

He, Alex and Robert Fay. 2022. *Digital Governance in China: Data, AI and Emerging Technologies, and Digital Trade*. Conference Report – Virtual Workshop, Centre for International Governance Innovation, Waterloo, ON, November 28. www.cigionline.org/publications/digital-governance-in-china-data-ai-and-emerging-technologies-and-digital-trade/.

Jaillant, Lise. 2022. *Archives, Access and Artificial Intelligence: Working with Born-Digital and Digitized Archival Collections*. Bielefeld, Germany: Bielefeld University Press. https://doi.org/10.1515/9783839455845.

Klosek, Katherine and Marjory S. Blumenthal. 2024. "Training Generative AI Models on Copyrighted Works Is Fair Use." *ARL Views* (blog), January 23. www.arl.org/blog/training-generative-ai-models-on-copyrighted-works-is-fair-use/.

Milmo, Dan. 2024. "'Impossible' to create AI tools like ChatGPT without copyrighted material, OpenAI says." *The Guardian*, April 8. www.theguardian.com/technology/2024/jan/08/ai-tools-chatgpt-copyrighted-material-openai.

Ministry of Government and Consumer Services. 2018. *Freedom of Information and Protection of Privacy Manual*. Toronto, ON: Government of Ontario. https://files.ontario.ca/books/mgcs-foi-privacy-manual-en-2021-09-02.pdf.

Moreton, Álvaro and Ariadna Jaramillo. 2021. "The problem of complete, irreversible data anonymization." AIhub, July 6. https://aihub.org/2021/07/06/the-problem-of-complete-irreversible-data-anonymisation/.

Paprica, Alison. 2023. "Exothermic innovation: Look beyond start-up costs and focus on the energy needed to keep change progressing." *Medium*, August 28. https://medium.com/@papricaalison/exothermic-innovation-407cab3c61c4.

Rustad, Arne. 2024. "Adaption of Generative Methods for Anonymization will Revolutionize Data Sharing and Privacy." *Medium*, January 17. https://towardsdatascience.com/adaption-of-generative-methods-for-anonymization-will-revolutionize-data-sharing-and-privacy-d35b6fe704a2.

SCC. 2021. *Canadian Data Governance Standardization Roadmap*. SCC, June 28. https://scc-ccn.ca/resources/publications/canadian-data-governance-standardization-roadmap.

Solow-Niederman, Alicia. 2024. "Can AI Standards Have Politics?" *UCLA Law Review* 71: 230–45.

Thompson, Graham. 2023. "Data Anonymization in AI: A Path Towards Ethical Machine Learning." *Security & Privacy* (blog), November 22. www.privacydynamics.io/post/data-anonymization-in-ai-a-path-towards-ethical-machine-learning/.

Vázquez, Favio. 2018. "Ontology and Data Science: How the study of what there is can help us be better data scientists." *Medium*, December 7. https://towardsdatascience.com/ontology-and-data-science-45e916288cc5.

Villalobos, Pablo, Anson Ho, Jaime Sevilla, Tamay Besiroglu, Lennart Heim and Marius Hobbhahn. 2022. "Will we run out of data? Limits of LLM scaling based on human-generated data." *arXiv*, June 4. https://doi.org/10.48550/arxiv.2211.04325.

Vorvoreanu, Mihaela, Amy Heger, Samir Passi, Shipi Dhanorkar, Zoe Kahn and Ruotong Wang. 2023. *Responsible AI Maturity Model: Mapping Your Organization's Goals on the Path to Responsible AI*. Microsoft, May 17. www.microsoft.com/en-us/research/publication/responsible-ai-maturity-model/.

Wilkinson, Mark D., Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg et al. 2016. "The FAIR Guiding Principles for scientific data management and stewardship." *Scientific Data* 3, 160018. https://doi.org/10.1038/sdata.2016.18.