

# Registry Data Storage and Curation

Megan Johnston<sup>1,7</sup>, Craig Campbell<sup>2</sup>, Rachel Hayward<sup>3</sup>, Mark Lowerison<sup>1,7</sup>,  
Vanessa K. Noonan<sup>4</sup>, Ted Pfister<sup>1,7</sup>, Colleen Maxwell<sup>5</sup>, Claire Marie Fortin<sup>6</sup>,  
Eric E. Smith<sup>1,7</sup>, Jean K. Mah<sup>1,7</sup>, Moira K. Kapral<sup>8</sup>, Nathalie Jette<sup>1,7,9</sup>,  
Tamara Pringsheim<sup>1,7</sup>, Lawrence Korngut<sup>1,7</sup>

Can J Neurol Sci. 2013; 40: Suppl. 2 - S35-S40

The storage of patient and medical information in a disease registry is a critical concept for consideration during registry design and development. The choice of data storage methods may influence the ability to access data in the future; the ability to store data long-term; and the ability to exchange data with other registries or research projects as required. Additionally choosing a data storage method involves a certain degree of uncertainty in an era that has gone from the file cabinet to the five inch floppy to the cloud in a matter of 35 years. In preparing this section of the guideline we reviewed available scholarly and grey literature resources; consulted with disease, registry, legal, ethics, privacy, and information technology (IT) experts; and consulted appropriate legislation and policy documentation in Canada.

## RELEVANT LITERATURE

Unfortunately our efforts to examine relevant literature in this topic area were unsuccessful. While there is a large body of IT literature on topics that may apply here, very little is specific to the Canadian context or the disease registry context. Where general principles applied we have reflected this as much as possible. Additionally, in some registry literature where mention to the issues of Data Storage and Curation were made we have noted this.

## Policy and Legislation

Many Canadian provinces and territories have specific legislation components that address information technology applications and criteria that must be met by applications collecting health information. As a result, disease registry projects need to consider their relevant legislation within the jurisdiction in which the database itself will be housed, and any other additional needs that could be demanded of the registry based on the other jurisdictions in which it operates. Table 5 features a list of relevant documentation by province.

When examining software products to determine the best fit for a registry application; evaluate the product specifications to ensure that all legislative requirements can be met. Table 6

outlines some of the common requirements for neurological registries in Canada and some software products available in 2012 that meet some or all of the requirements.

## OTHER CONSIDERATIONS

### Storage Considerations

The type of database selected for a disease registry project will depend on a number of factors determined early in the registry development including: the expected number of records (database size); the expected number of users (database clients); the expected duration of the registry (length of data storage); the type of data being stored (data type); and the duration of the data storage after the registry project is complete. For example, in Canada, clinical trial data are required to be stored for 25 years under Part C Division 5 of the Food and Drug Regulations [C.05.012], however little consideration is typically given to the format of the storage of clinical trial data and whether or not this will remain accessible 25 years in the future. With electronic data storage, such considerations must not be underestimated. If registries are capturing both observational and clinical trial data, there may also be a need to store the observational data much longer than might normally be the case or to have the registry modules separated so that data from clinical trials can be stored for the longer time frame. These considerations should be made in advance of registry set up as they may impact the type of consent provided by patients in the area of data storage.

In addition to the above considerations disease registry projects may also want to consider Canadian legislation and privacy considerations with respect to data storage location. The following aspects should be considered:

- A) Server Model (e.g. single server, dual server, or cloud server/storage?)
- B) Physical Location of Servers (e.g. country, province, institution)
- C) Physical Server Access (e.g. controlled, secure?)
- D) Network location of Server (e.g. secured, visible, access controls)

From the <sup>1</sup>University of Calgary, Calgary, Alberta; <sup>2</sup>Western University – London Health Sciences Centre, London, Ontario; <sup>3</sup>Alberta Office of the Information and Privacy Commissioner, Edmonton, Alberta; <sup>4</sup>Rick Hansen Institute, Vancouver, British Columbia; <sup>5</sup>University of Waterloo, Waterloo, Ontario; <sup>6</sup>Canadian Institutes of Health Information, Toronto, Ontario; <sup>7</sup>Hotchkiss Brain Institute, University of Calgary, Calgary, Alberta; <sup>8</sup>University of Toronto, Toronto, Ontario; <sup>9</sup>Institute for Public Health, University of Calgary, Calgary, Alberta.

FINAL REVISIONS SUBMITTED JANUARY 28, 2013.

Correspondence to: Lawrence Korngut, Clinical Neurosciences, 480060, 4th Floor Administration, South Health Campus, 4448 Front Street SE, Calgary, Alberta, T3M 1M4, Canada. Email: Lawrence.korngut@albertahealthservices.ca.

- E) Database user access (e.g. data access permission levels; authentication mechanisms)
- F) Hardware and Software security controls (e.g. firewalls, encryption)

**Database Genres**

When selecting a database genre (database type) consider the complexity of the data processing that will be required during registry operation and the organizational resources available for the management of the database. Table 7, adapted from Brian Westrich, University of Minnesota<sup>151</sup> may be a useful tool during these considerations.

Following identification of the required database genre it will be necessary to select a specific software product with which to

execute the database. Considerations during this process will include organizational assets (e.g. institutional licenses or IT services); budget (consider using open source software products if budget is small) and the development timeline. Additionally considerations must be made regarding the software product's ability to meet data storage requirements associated with legislation (See Table 5). Finally a key consideration during this stage involves the database size. The larger and more complex the database, the more important it becomes to select a software product that can create an efficient and readily accessible database while optimizing storage space. To this end, one must consider the structure of the database created by each database genre product. Additionally storage space will be impacted by the format of the data stored in the database.

**Table 5: Relevant Legislation and Policy Relating to Software Considerations**

Province/Territory	Best Practice/Guidelines Document
Alberta	Personal Information Protection Act (PIPA) PIPA Advisory #8 Implementing Reasonable Safeguards ( <a href="http://www.oipe.ab.ca/Content_Files/Files/Publications/PIPA_Advisory_8_Reasonable_Safeguards2007.pdf">http://www.oipe.ab.ca/Content_Files/Files/Publications/PIPA_Advisory_8_Reasonable_Safeguards2007.pdf</a> ) <sup>(111)</sup>
	Alberta Electronic Health Record Regulation ( <a href="http://www.qp.alberta.ca/documents/Regs/2010_118.pdf">http://www.qp.alberta.ca/documents/Regs/2010_118.pdf</a> ) <sup>(112)</sup>
	FOIP Guidelines and Practices Chapter 8. Records and Information Management ( <a href="http://www.servicealberta.ca/foip/documents/chapter8.pdf">http://www.servicealberta.ca/foip/documents/chapter8.pdf</a> ) <sup>(113)</sup>
	Developing Records Retention and Disposition Schedules ( <a href="http://www.rimp.gov.ab.ca/publications/pdf/SchedulingGuide.pdf">http://www.rimp.gov.ab.ca/publications/pdf/SchedulingGuide.pdf</a> ) <sup>(114)</sup>
	Health Information Act Guidelines and Practices Manual ( <a href="http://www.health.alberta.ca/documents/hia-guidelines-practices-manual.pdf">http://www.health.alberta.ca/documents/hia-guidelines-practices-manual.pdf</a> ) <sup>(115)</sup>
	FOIP Guidelines and Practices ( <a href="http://www.servicealberta.ca/foip/resources/guidelines-and-practices.cfm">http://www.servicealberta.ca/foip/resources/guidelines-and-practices.cfm</a> ) <sup>(116)</sup>
British Columbia	Physicians & Security of Personal Information ( <a href="http://www.oipe.bc.ca/tools-guidance/guidance-documents.aspx">http://www.oipe.bc.ca/tools-guidance/guidance-documents.aspx</a> ) <sup>(117)</sup>
	Information Management and Information Technology Management ( <a href="http://www.fin.gov.bc.ca/ocg/fmb/manuals/CPM/12_Info_Mgmt_and_Info_Tech.htm">http://www.fin.gov.bc.ca/ocg/fmb/manuals/CPM/12_Info_Mgmt_and_Info_Tech.htm</a> ) <sup>(118)</sup>
	FOIPP Act Policy and Procedures Manual ( <a href="http://www.cio.gov.bc.ca/cio/priv_leg/manual/sec30_39/sec30_page?">http://www.cio.gov.bc.ca/cio/priv_leg/manual/sec30_39/sec30_page?</a> ) <sup>(119)</sup>
	Information Security Policy ( <a href="http://www.cio.gov.bc.ca/local/cio/informationsecurity/policy/isp.pdf">http://www.cio.gov.bc.ca/local/cio/informationsecurity/policy/isp.pdf</a> ) <sup>(120)</sup>
Manitoba	University of Manitoba Safe Computing Topics ( <a href="http://www.oit.umn.edu/safe-computing/topics/index.htm">http://www.oit.umn.edu/safe-computing/topics/index.htm</a> ) <sup>(121)</sup>
	Respecting Privacy: A Compliance Review Tool for Manitoba's Information Privacy Laws: A Special Report <a href="http://www.ombudsman.mb.ca/pdf/Special%20Report%20English%20CRT%20-%20Oct%207.pdf">http://www.ombudsman.mb.ca/pdf/Special%20Report%20English%20CRT%20-%20Oct%207.pdf</a> <sup>(122)</sup>
New Brunswick	Guidelines for Custodians to assess compliance with the Personal Health Information Privacy and Access Act (PHIPPA) ( <a href="http://www.gnb.ca/0051/acts/pdf/7133%20E2%82%AC%20English%20long%20list%203s.pdf">http://www.gnb.ca/0051/acts/pdf/7133%20E2%82%AC%20English%20long%20list%203s.pdf</a> ) <sup>(123)</sup>
Newfoundland	The Personal Health Information Act (Resources) ( <a href="http://www.health.gov.nl.ca/health/PHIA/">http://www.health.gov.nl.ca/health/PHIA/</a> ) <sup>(124)</sup>
Nova Scotia	Personal Health Information Legislation for Nova Scotia ( <a href="http://novascotia.ca/dhw/phia/custodians.asp">http://novascotia.ca/dhw/phia/custodians.asp</a> ) <sup>(125)</sup>
	Privacy Impact Assessment Template ( <a href="http://www.gov.ns.ca/just/IAP/docs/Appendix%20B%20PIA%20Template.pdf">http://www.gov.ns.ca/just/IAP/docs/Appendix%20B%20PIA%20Template.pdf</a> ) <sup>(126)</sup>
Nunavut	
Ontario	IPC Ontario Privacy and Confidentiality When Working Outside the Office ( <a href="http://www.ipc.on.ca/images/Resources/up-num_20.pdf">http://www.ipc.on.ca/images/Resources/up-num_20.pdf</a> ) <sup>(127)</sup>
	Manual for the Review and Approval of Prescribed Persons and Prescribed Entities ( <a href="http://www.ipc.on.ca/images/Findings/process.pdf">http://www.ipc.on.ca/images/Findings/process.pdf</a> ) <sup>(128)</sup>
Prince Edward Island	According to the Forms and Resource Materials section of the Information and Privacy Commissioner's website ( <a href="http://www.assembly.pe.ca/index.php3?number=1013951">http://www.assembly.pe.ca/index.php3?number=1013951</a> ) <sup>(129)</sup> PEI's FOIP Act is based on Alberta's FOIP Act and cites the Guidelines and Practices: 2009 Edition ( <a href="http://www.servicealberta.ca/foip/resources/guidelines-and-practices.cfm">http://www.servicealberta.ca/foip/resources/guidelines-and-practices.cfm</a> ) <sup>(116)</sup> as a useful reference. Chapter 9 lists technical safeguards.
Quebec	Minimum Requirements for the Security of Computerized Records of Health and Social Services Network Clients ( <a href="http://www.cai.gouv.qc.ca/documents/CAI_G_securite_doss_info_rsss_eng.pdf">http://www.cai.gouv.qc.ca/documents/CAI_G_securite_doss_info_rsss_eng.pdf</a> ) <sup>(130)</sup>
	Exigences minimales relatives à la sécurité des dossiers informatisés des usagers du réseau de la Santé et des Services sociaux <a href="http://www.cai.gouv.qc.ca/documents/CAI_G_securite_doss_info_rsss.pdf">http://www.cai.gouv.qc.ca/documents/CAI_G_securite_doss_info_rsss.pdf</a> <sup>(131)</sup>
	CAI Quebec ( <a href="http://www.cai.gouv.qc.ca/english/">http://www.cai.gouv.qc.ca/english/</a> ) ( <a href="http://www.cai.gouv.qc.ca/">http://www.cai.gouv.qc.ca/</a> )
Saskatchewan	Saskatchewan Archives Board Records Management Policies and Guidelines ( <a href="http://www.saskarchives.com/services-government/record-management-policy-and-guidelines">http://www.saskarchives.com/services-government/record-management-policy-and-guidelines</a> ) <sup>(132)</sup>
	Security Controls for Protection of Personal Information ( <a href="http://www.justice.gov.sk.ca/TTOSecurityControlsforProtectionofPersonalInformation.pdf">http://www.justice.gov.sk.ca/TTOSecurityControlsforProtectionofPersonalInformation.pdf</a> ) <sup>(133)</sup>
	Government of Saskatchewan Resources and Tools Security ( <a href="http://www.justice.gov.sk.ca/AP_Security">http://www.justice.gov.sk.ca/AP_Security</a> ) <sup>(134)</sup>
Northwest Territories	GNWT INFORMATION TECHNOLOGY Electronic Information Security ( <a href="http://www.fin.gov.nt.ca/documents/ocio/ppse/6003.00.27%20-%20Standards%20-%20Electronic%20Information%20Security.pdf">http://www.fin.gov.nt.ca/documents/ocio/ppse/6003.00.27%20-%20Standards%20-%20Electronic%20Information%20Security.pdf</a> ) <sup>(135)</sup>
Yukon	ATIPP Compliance Assessment ( <a href="http://www.ombudsman.yk.ca/uploads/general/ACA_ATIPP_Compliance_Assessment_August_2011.pdf">http://www.ombudsman.yk.ca/uploads/general/ACA_ATIPP_Compliance_Assessment_August_2011.pdf</a> ) <sup>(136)</sup>
	Yukon Information and Privacy Commissioner Privacy Breach Checklist ( <a href="http://www.ombudsman.yk.ca/uploads/general/ATIPP_Privacy_Breach_Checklist_2011.pdf">http://www.ombudsman.yk.ca/uploads/general/ATIPP_Privacy_Breach_Checklist_2011.pdf</a> ) <sup>(137)</sup>
	Yukon Information and Privacy Commissioner Best Practice: Responding to a Privacy Breach ( <a href="http://www.ombudsman.yk.ca/uploads/general/ATIPP_Best_Practice_Privacy_Breach_Response.pdf">http://www.ombudsman.yk.ca/uploads/general/ATIPP_Best_Practice_Privacy_Breach_Response.pdf</a> ) <sup>(138)</sup>
	Privacy Impact Assessment ( <a href="http://www.ombudsman.yk.ca/uploads/general/PRIVACY%20IMPACT%20ASSESSMENT.pdf">http://www.ombudsman.yk.ca/uploads/general/PRIVACY%20IMPACT%20ASSESSMENT.pdf</a> ) <sup>(139)</sup>
Canada	Operational Security Standard: Management of Information Technology Security (MITS) ( <a href="http://www.tbs-sct.gc.ca/pol/doc-eng.aspx?id=12328&amp;section=text">http://www.tbs-sct.gc.ca/pol/doc-eng.aspx?id=12328&amp;section=text</a> ) <sup>(140)</sup>
	Guidance Document: Taking Privacy into Account Before Making Contracting Decisions ( <a href="http://www.tbs-sct.gc.ca/atip-airtp/tpa-ppc/tpa-ppc06-eng.asp">http://www.tbs-sct.gc.ca/atip-airtp/tpa-ppc/tpa-ppc06-eng.asp</a> ) <sup>(141)</sup>
	Electronic Health Record (EHR) Privacy and Security Requirements Reviewed with Jurisdictions and Providers ( <a href="https://knowledge.inforoute.ca/EHRSRA/doc/EHR-Privacy-Security-Requirements.pdf">https://knowledge.inforoute.ca/EHRSRA/doc/EHR-Privacy-Security-Requirements.pdf</a> ) <sup>(142)</sup>
	Health Canada Final Audit Report – Audit of Information Technology (IT) Security ( <a href="http://www.hc-sc.gc.ca/ahc-asc/pubs/audit-verif/2011-04/index-eng.php#_Toc2008">http://www.hc-sc.gc.ca/ahc-asc/pubs/audit-verif/2011-04/index-eng.php#_Toc2008</a> ) <sup>(143)</sup>
Other	ISO/IEC 27002:2005 Information technology – Security techniques – Code of practice for information security management ( <a href="http://www.iso27001security.com/html/27002.html#">http://www.iso27001security.com/html/27002.html#</a> ) <sup>(144)</sup>

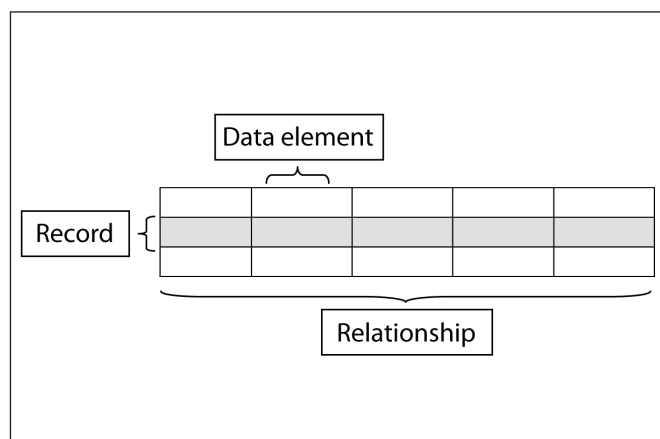
**Table 6: Software features by product**

FEATURE	Open-Source	Local Server Install	Authentication/Platform	Password Controlled User Level-Authentication	User Action Log	Data encryption	Interacts with Third-Party Data Sources	Workflow Management	Patient Portal
<b>PRODUCT</b>									
<b>REDCap</b> (Produced by Vanderbilt University) <a href="http://www.project-redcap.org">www.project-redcap.org</a> <sup>(146)</sup>	✓	✓ (if you are a participating site)	LDAP, Shibboleth or Table-Based (user selected)	✓	✓	Can be installed using additional software products	✓		
<b>ClinicalPursuit</b> <a href="http://www.patientregistrysoftware.com">www.patientregistrysoftware.com</a> <sup>(147)</sup>		optional	Microsoft.NET	✓	✓		✓		
<b>i2b2 – Informatics for Integrating Biology and the Bedside</b> <a href="http://www.i2b2.org">www.i2b2.org</a> <sup>(148)</sup>	✓	✓	AJAX	✓	✓		✓	✓	
<b>Patient Crossroads</b> <a href="http://www.patientcrossroads.com">www.patientcrossroads.com</a> <sup>(149)</sup>				✓					✓
<b>Axiom Clarinet</b> <a href="http://www.certus-tech.com">www.certus-tech.com</a> <sup>(150)</sup>		✓	J2EE	✓				✓	

**Database Structures**

The type of database structure that is selected will influence many factors impacting the operation of the registry. These factors might include: computer hardware infrastructure; registry stability and performance; data entry and recall speed; and reporting capability.

**Relational databases** – This type of database (see Figure 2) is still a very common format created by many software products on the market however it can come with some significant limitations if data sizes are large.<sup>152</sup> These databases store data in a defined record where the common location of the data elements contained within the record is the sole logic between the data elements within the record. This limits the granularity of the database to the record level (i.e. data cannot be examined within a record except if the full record is recalled). As a result processing time to read and write records is high; total disk storage required for the database is high; and modifications to records require the whole record to be rewritten.



**Figure 2: Relational Database Structure**

**Table 7: Database Genre Decision Support Tool**

Complexity	Minimum required data resources	Database genre
Storage only	Paper case report forms	Non-automated file storage (paper based and/or electronic copies of paper forms).
Electronic storage	Computer	Word processor or basic file storage software. Basic file backup to external media (CD-ROM or DVD).
Structuring – Data that is stored needs to have different “fields” or “pieces”.	Semi-skilled staff	Spreadsheet. Note that this is still a storage only task and no analysis is required.
Relating – Data is stored in fields and there is a need to define relationships or examine relationships between the fields.	Computer staff (part-time)	Personal database tools (e.g. Access). These tools feature simple data form and query design tools. Multiple data tables can be created and relationships between them can be defined. Analysis required is simplistic.
Complex, high volume – There will be large amounts of data between which complex relationships exist. This may also involve the need to have simultaneous access by multiple users.	Computer staff (full-time)	Industrial database tools (e.g. Oracle, mySQL). These tools allow for all of the features of Personal database tools but also allow for logging of user transactions; simultaneous access and updates by multiple users and complex query construction (for example, construction of data sub-sets). These software products may also allow database architecture to span over multiple servers for operation and storage.
Highly specific or specialized – The type of data being collected; the data collection process and/or the queries and analysis required of the data require customization beyond that available in standard tools.	Highly skilled computer programmers. High performance computing equipment. This type of solution may also require custom networking.	Programming languages (e.g. Java, C-plus). These tools may operate in conjunction with industrial database tools or other library structures to fully enable the required database architecture.

**Columnar databases** – This type of database (see Figure 3 below) is increasingly adopted due to the increase in analytical simplicity found through this method when compared to relational databases. These databases store information by column with all values within a column being stored as a single dataset (i.e. these datasets are made up of data from multiple “records”).<sup>152</sup> A key advantage of this format is that “parts of records” from a relational database perspective can be analyzed and written or rewritten. This feature increases the speed with which data processing can be accomplished. However, the trade off here is that recalling records requires the assembly of data values across multiple columns into a pre-determined format which if the number of columns is large (complex dataset) or the number of requests is large (many users) may impact database performance.

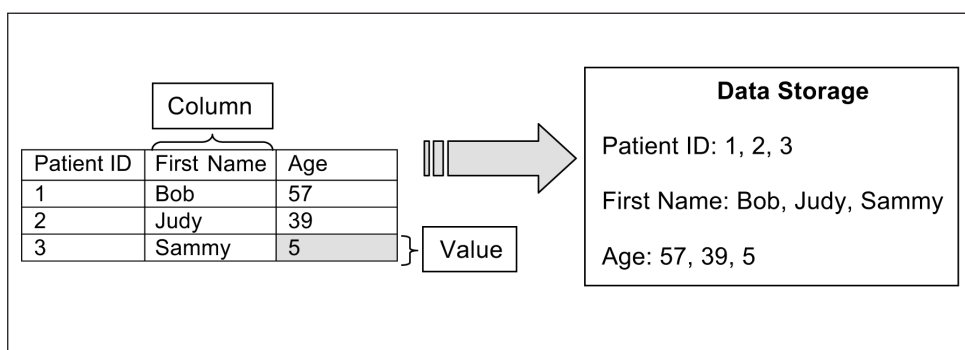


Figure 3: Columnar Database Structure

**Correlation databases** – This type of database (see Figure 4 below) stores each data value only once and then stores references allowing collocation of appropriate values for each “record” using descriptive metadata. Like a card catalogue, metadata stores information on what values are required for each “record” and where each value can be found allowing programming that reads the metadata to reassemble each “record” when required. These databases have similar advantages to columnar databases in terms of partial record access and writing actions. However due to the low storage volume of correlation databases, their performance often exceeds columnar databases.<sup>152</sup>

Once the database genre and database structure are selected, the final considerations are the database formatting and the configuration of adequate backup infrastructure.

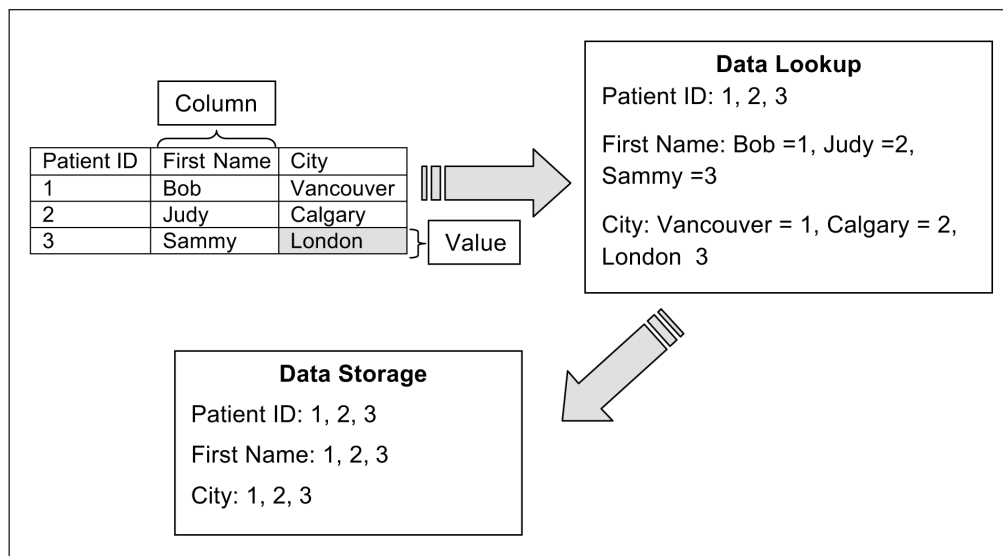


Figure 4: Correlation Database Structure

**Data Formats**

There are three ways to physically store digital data available on the current market<sup>153</sup>:

- 1) Magnetic storage (e.g. magnetic tape or hard drive)
- 2) Solid state (e.g. flash memory)
- 3) Optical (e.g. Blu-Ray, DVD, CD-ROM)

Table 8 discusses some of the considerations for each method. Clearly demonstrated by the information in Table 8, the choice of backup infrastructure is best addressed by choosing multiple physical storage formats. It should be a regular practice to perform backups at pre-determined intervals to a hard drive space and then periodic backups to an optical format.

In addition to physical storage format, in registry development one must consider file format obsolescence. This is the state during which the digital format of the file is no longer readable due to changes in technology and file formatting practices. File format obsolescence is independent of physical storage format and is to do with the actual digital format of the files on the physical storage format. Both are important considerations when storing data long term.

As it is impossible to define the file formats of the future and indeed to define file formats that will not go obsolete this guideline instead recommends a risk assessment approach to addressing this concern. This risk assessment approach is derived from File Format Obsolescence Risk Decision Support System (Version 1.1 released November 2007).<sup>155</sup>

**Table 8: Data Format Considerations**

Data Format	Life Expectancy	Pros	Cons
Magnetic tape	0.7 – 1083 years depending on storage temperature, humidity, availability of error-correction coding <sup>(153)</sup>	<ul style="list-style-type: none"> <li>• Readily available</li> <li>• Cheap</li> <li>• Convenient</li> <li>• Sizeable</li> </ul>	<ul style="list-style-type: none"> <li>• Oldest technology</li> <li>• Many known impacts on life expectancy</li> <li>• Reliability depends on manufacturing</li> </ul>
Hard disk drives	Limited knowledge available but may range from as early as 3 months independent of utilization <sup>(154)</sup>	<ul style="list-style-type: none"> <li>• Readily available</li> <li>• Convenient</li> <li>• Sizeable</li> </ul>	<ul style="list-style-type: none"> <li>• Limited capacity (currently in terabyte (TB) range)</li> <li>• Failure is typically catastrophic</li> </ul>
Flash memory (solid state drives or memory sticks)	10 – 13 years without use. May extend up to 100 years with active management. <sup>(153)</sup>	<ul style="list-style-type: none"> <li>• Readily available</li> <li>• Convenient</li> <li>• Small/Portable</li> </ul>	<ul style="list-style-type: none"> <li>• Current size limitation is 8 GB per unit.</li> <li>• Loss is inevitable without active management.</li> </ul>
Optical media (ROM)	20 – 12,000 years. <sup>(153)</sup>	<ul style="list-style-type: none"> <li>• Readily available</li> <li>• Requires little technical knowledge</li> <li>• Potentially lengthy life expectancy.</li> <li>• Permanent (once written it is only readable).</li> <li>• Multiple densities available</li> </ul>	<ul style="list-style-type: none"> <li>• Requires dedicated drive technology for reading and writing.</li> <li>• Drive technology may become obsolete.</li> <li>• Need for secure physical storage location in which to retain media.</li> <li>• Limited unit storage size.</li> </ul>
Optical media (recordable)	Light and heat dependent but can be as low as a few hours in direct sunlight <sup>(154)</sup>	<ul style="list-style-type: none"> <li>• Readily available</li> <li>• Requires little technical knowledge</li> <li>• Multiple formats available</li> <li>• Rewritable (non-permanent)</li> </ul>	<ul style="list-style-type: none"> <li>• Need for secure and dark physical storage location in which to retain media.</li> <li>• Requires dedicated drive technology for reading and writing.</li> <li>• Drive technology may become obsolete.</li> <li>• Limited unit storage size.</li> </ul>

- 1) Is the file format a standard base coding format? (e.g. UNICODE, ASCII)
  - a. Yes – This file format is low risk.
  - b. No – Continue to Question 2.
- 2) Is the file format referenced in any searchable information resources?
- 3) Is there a known support end date for the file format?
  - a. Yes – How many years until the support end date? (if within your long term storage needs, consider an alternate format).
- 4) How many years since this file format version was released?
  - a. Are new versions available or on the horizon?
- 5) What is the primary rendering software needed for this file format? (i.e. what program is needed to read the files)
  - a. Identified – is this software available to you?
- 6) Does the primary rendering software have critical hardware or software needs that might not be available in the future?
  - a. Yes – what equipment/software is required and is it available to you?
- 7) Are there alternate rendering software solutions available? If so, what are they and what are their hardware/software requirements?
  - a. Identified – how many of these are available to you with all their requirements?
- 8) Are there other means of providing safe and effective access? (e.g. custom coding, open source applications?)
- 9) What is the total number of access methods available for your chosen file format? Include all primary and alternate methods.
  - a. 0 – Extremely high risk, consider your data as lost.
  - b. 1 – High risk, consider alternate formats if possible
  - c. 2 to 5 – Medium risk, ensure hardware/software requirements for access are documented and retained if possible.
  - d. 5 or more – Low risk, you can proceed with implementation.

### Data Migration

Data migration involves the process of moving data from one source to another where the structure of the data will change.<sup>5</sup> Data migration might be necessary if a registry platform becomes obsolete (either due to changes in software design or due to software discontinuation); if software cost becomes an issue (in the case of proprietary software platforms especially); or if additional functionality not planned in the initial registry design is required. Further data migration may be required if physical server characteristics or locations change; if data ownership requirements or personnel change; or to meet larger IT infrastructure expectations within the host organization. While it is possible to do data migration manually, the time investment will be considerable and efforts should be made to select software that can mediate an automated migration. For an automated migration to be successful and detailed data map correlating every data type from the old system into the new system must be created.<sup>5</sup> A plan for handling inconsistent data should be created at the outset and revised if any additional issues are raised during the migration process.<sup>5</sup> Following migration, quality assurance activities should be conducted to ensure that the data in the new system has been transferred as expected. Overall, a simple project management methodology made popular by W. Edwards Deming of “Plan-Do-Check-Act” (the Shewhart Cycle) can be a great approach for a data migration project. First, plan the data migration including required staffing and software/hardware resources; project timelines and any server down time that may occur. As previously mentioned, ensure that this plan features a detailed map of the data migration from the old system to the new system. Next perform a small test migration. Following the test, enact your quality assurance plan and evaluate if the desired results have been achieved. Once you have reviewed the test results, either revisit the plans and retest, or proceed with the full migration.

**Data Curation**

Data curation is defined by Lord and MacDonald as “the activity of managing and promoting the use of data from its point of creation, to ensure it is fit for contemporary purpose, and available for discovery and reuse”.<sup>156</sup> This activity may include simple data management activities, enriching or adding value to data, the sharing of data, and the preservation of data for a later use. Data curation is a critical activity for the creation and maintenance of successful disease registries. While this guideline cannot define an individual registry’s curation plan, below are some key points to ensure creation of a complete data management plan. These key points are taken from a Data Management Plan checklist produced by the Digital Curation Centre.<sup>157</sup>

- 1) Data types – understand what types of data will be collected in the registry. Ensure that data are defined using a data dictionary.
- 2) Data formats – consider the format of each type of data (e.g. text, alphanumeric, date etc). List all the possible formats that will be collected across the registry.
- 3) Standards – likely partly outlined in the data dictionary but ensure that there is clear definition of what data will be accepted and rejected. Additionally document if data will be compared to other sources and any associated standards dictated by this relationship.
- 4) Capture methods – document all methods of data capture and data flow into the data repository (database).
- 5) Data output – consider what content is being created by the project and document this. For example does the project simply produce raw data for further use or does the project produce derived data.<sup>158</sup>
- 6) Storage - what storage space is required for the data output? See Storage Considerations in this Guideline for more information.
- 7) File formats – what file formats will be used and why? Ensure that you document your analysis of file format considerations and risks in the data management plan. See Storage Considerations in this Guideline for more information.
- 8) Future uses – consider what the future uses or reuse of the data output and/or original data might be. What will be required to ensure these future uses/reuse can occur.
- 9) Sharing – ensure that consideration of whom might share the data and all associated ethical, legal and logistical issues are outlined and addressed.
- 10) Access – who will have access to the data and what are the access controls?
- 11) Existing data – are there existing data that are required or beneficial to the project. What constraints or considerations are present as a result? Is new data production actually needed? What is the value of the new data? What access is required to obtain existing datasets?
- 12) Data quality – what is your plan for data quality assurance and control?
- 13) Documentation – what documentation is required to ensure that data make sense in isolation? Consider that the context required may be stored with the data itself using metadata.
- 14) Metadata – if metadata will be included ensure you have considered how they will be created, maintained and stored.

- 15) Intellectual Property – ensure that the ethical and legal considerations associated with existing data and new data have been considered and addressed. See the Ethics & Privacy section of this guideline for more information.
- 16) Accountability – who is responsible for the data and who are the delegates of this authority if applicable? How is accountability assigned (e.g. legislation; institutional policy)? How will accountability be transitioned if required?

**Data Management Plan**

A data curation document will be part of a larger data management plan. The data management plan will include additional aspects such as:

- Who manages the data?
- Where, how and when will data be backed up?
- What mechanisms are in place for error tracking and change logging?
- Who is responsible for addressing changes, errors and trouble? What is the process for addressing changes, errors and trouble?
- What security systems are in place to protect the data?
- What is the process if there is a security breach?

For assistance creating a comprehensive data management plan, consider utilizing the DMPTool found at <https://dmp.cdlib.org/>.<sup>159</sup>

**RECOMMENDATIONS**

- ✓ In the context of applicable Canadian legislation consider the following items with respect to data storage:
  - o Server Model
  - o Physical Location of Servers and Access
  - o Network Location of Servers
  - o User Access levels and permissions
  - o Hardware and software security controls
- ✓ Consider the complexity of your storage needs and the required personnel and software resources to maintain them.
  - o Maximize organizational assets such as existing software licenses or discounts.
  - o Wherever possible utilize open source software to minimize development and ongoing costs.
  - o Document and plan your development timeline.
  - o Ensure you have planned for adequate storage space and database size/functionality. Assess required computer hardware to facilitate desired access times; registry stability and needed reporting capabilities.
- ✓ Choose multiple data storage formats for short and long-term data backup. Ensure backup plans meet necessary legislation and policy expectations. Document data backup procedures and schedule in the data management plan.
- ✓ Assess file format storage risk.
- ✓ Create data curation and data management plans.