# BUILDING A SCALEABLE NATIONAL REPOSITORY INFRASTRUCTURE FOR CANADA

Todd Trann, University of Saskatchewan
Alex Garnett, Simon Fraser University

CANHEIT June, 2016

# What Is A Data Repository?

Main functions: **digital storage, retrieval and preservation** of research data

**Institutional Repositories** are run by colleges and universities

- cIRcle (Run by UBC)

**Regional Repositories** are run by groups of universities

- Scholars Portal (Run by OCUL)

**Domain Specific Repositories** are run by organizations or communities of researchers

- BioGRID (Biology; run by the BioGRID team)

There are thousands of research repositories in the world (see: re3data.org)

# YARR

- Current repositories don't scale to very large data sets

- Hard to discover data sets in existing repositories

- Some researchers don't have access to a suitable repository

# Project Overview

Designing infrastructure for a national research data repository
- Partnerships between Compute Canada, Globus and CARL (Canadian Association of Research Libraries)
- Using Compute Canada hardware

Designing discovery platform
- Indexing metadata of existing repositories
- Sharing back a search API
- Improving cross-repository discovery
- Facilitating domain specific search development

# Planned Features

- **Federated Storage Model:** Storage and repositories can be distributed, owned and operated by organizations / institutions

- **Scalable:** Many files; large files; large total amounts of storage

- **National Data Discovery:** Single search to discover data, regardless of location

# Planned Features (cont'd)

- **Data Preservation:** Use Archivematica to automatically preserve data

- **Suitable for broad range of data types**

- **Automatic geographic data replication:** We intend to leverage Compute Canada national data infrastructure being deployed in 2016

# System Diagram



**Search**

**Submit**

Globus Publication

Archive-matica

Globus Connect

CC Storage

**Compute Canada**

Globus Connect

Institutional Storage

**Institution**

Globus Search Webapp

Globus Search Engine

**Globus**

*Metadata Harvesting*

Regional Repository

Institutional Repository

Domain Repository

**Repositories**

# Metadata Overlay

Search:
victoria
title:victoria
author:victoria
gmd.contentInfo:victoria

Globus Search Webapp

Globus Search Engine

OAI-PMH Metadata Harvesting

Regional Repository

dublin core
datacite

Institutional Repository

dublin core

dublin core
gmd

Domain Repository

# Preservation Automation

# Discovery Landscape

- A University of British Columbia/CARL-Portage group (led by Eugene Barsky) is doing an environmental scan of data discovery systems, and considering options for the design of a coordinated, national discovery service for research data in Canada

- The white paper from this scan will inform finalization of the discovery interface

# Landing Page

# Bilingual Interface

# Submitting An Item

# Requesting Preservation

# Submission List



**ENDR NRDR**

Browse & Discover   **Data Publication Dashboard**   FR EN   Account ▾

## Your Submissions

Below are listed your previous submissions that have been accepted into the archive.

There are **8** datasets in the main archive that were submitted by you.

| Issue Date | Title | Author(s) |
|---|---|---|
| 8-Jul-2014 | ArcticNet 1305 - Northwest Passage CTD data | *Gratton, Yves* |
| 25-Sep-2014 | Acoustic data in the Beaufort Sea: BREA marine fishes 2012-2014 | *Fortier, Louis; Geoffroy, Maxime; Majewski, Andrew* |
| 17-Feb-2015 | Effect of movement between Salix arctica individuals on Gynaephora groenlandica caterpillar growth rates at Alexandra Fiord, Nunavut, Canada, 2013 | *Henry, Gregory; Greyson-Gaito, Christopher; Barbour, Barbour* |
| 15-Dec-2015 | Gastropod surveys in south-eastern Victoria Island | *Sullivan, Josh; McLennan, Donald; Kutz, Susan; Tomaselli, Matilde* |
| 16-Jul-2015 | Transects for pellet counts in tundra ecosystems around Cambridge Bay, Victoria Island | *McLennan, Donald; Anablak, Cathy* |
| 16-Dec-2015 | Bathymetry of a section of east end of Cambridge Bay, near the future site of the sea water intake for CHARS, Cambridge Bay, Nunavut | *McLennan, Donald* |
| 8-Jul-2015 | Identification and characterisation of subsurface flows in the Apex River, Iqaluit, Nunavut | *Franssen, Jan; Chiasson-Poirier, Gabriel; Lamoureux, Scott* |
| 15-Dec-2015 | Preliminary freshwater sampling of selected streams in the watershed of Greiner Lake, Cambridge Bay, Nunavut. | *McLennan, Donald* |

# Preservation Processing

# Starting A Search

# Search Results

# Item Page

# Item Page



AVOCADO RESEARCH EMAIL COLLECTION
hdl:11272/RXCHT
Version: 1 – Released: Thu Dec 03 14:00:55 PST 2015

**CATALOGING INFORMATION** | Data & Analysis | Comments (0) | Versions

If you use these data, please add the following citation to your scholarly references. Why cite?

**Data Citation**

Oard, Douglas; Webber, William; Kirsch, David; Golitsynskiy, Sergey, 2015-02-16, "Avocado Research Email Collection", http://hdl.handle.net/11272/RXCHT Linguistic Data Consortium [Distributor] V1 [Version]

**Citation Format** Print

**Data Citation Details**

| | |
|---|---|
| Title | Avocado Research Email Collection |
| Study Global ID | hdl:11272/RXCHT |
| Other ID | Linguistic Data Consortium: LDC2015T03; ISBN: 1-58563-704-1; ISLRN: 102-408-869-995-0 |
| Authors | Oard, Douglas; Webber, William; Kirsch, David; Golitsynskiy, Sergey |

# Dataset File Transfer



Globus Transfer
(large file capability)

# Further Information

Public web site:

www.computecanada.ca/rdm