

Canadian Research Knowledge Network
Réseau canadien de documentation pour la recherche

Héritage Collection Update: RG-10 ICR Pilot Project

Jason Friedman, Senior Manager, Heritage Services



crkn-rcdr.ca

1

Canadiana Collections: Content

Canadiana
Digitized monographs, serials, government publications, and maps

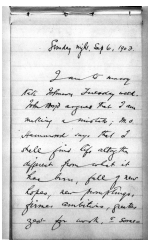
Héritage
Digitized microfilm of archival primary-source materials from Library and Archives Canada

2

Héritage collection

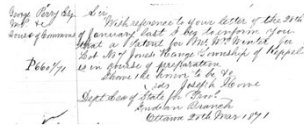
- Partnership with Library and Archives Canada (LAC)
- "C", "H" and "T" series of LAC microfilm
- 42+ million pages; 25,000+ reels, 900 collections
- Still in progress: 150,000 pages added annually



3

Discovery challenges

- Digitized microfilm
 - Limited finding aids/indexes
 - Limited metadata
 - Limited search capabilities
 - Lots of Handwriting means no OCR, no full text search






CRKN RCDR

4

ICR Pilot Project

5

Intelligent Character Recognition (ICR)

-  Uses AI to perform character recognition
-  Can be used on handwritten materials
-  Environmental scan conducted in Fall 2022

6

Transkribus

- Produced by READ-COOP
- Used by researchers/historians working with parts of the Héritage collection
- Initial sample sent for testing
 - 1-2% error rate for typewritten materials
 - 5-7% error rate for handwritten materials
- Next step: pilot project

7

ICR Pilot Project

- Partnership with LAC
- Full scope: 2 million images
- Began in March 2023
- Initial batch of 600K+
- Additional 500K identified
- Approximately 1 million more images to be identified

8

Why RG-10?

- High priority collection
- Heavily used
- Lots of handwritten content
- Support claims researchers
- Variety of document types/quality

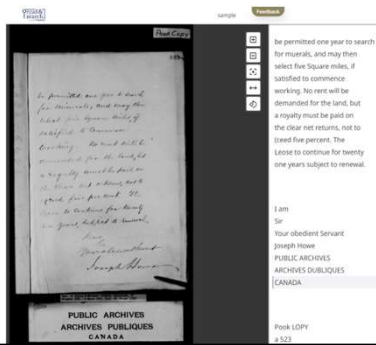
9

First Batch: Departmental Letterbooks

- All are handwritten with relatively consistent writing and format
- Covers a breadth of content of interest for variety of research topics
- They could benefit from ICR as they are poorly described
- Time period, 1871-1929 (low risk of sensitive content)
- Existing transcription data for some content from University of British Columbia/University of Saskatchewan for model training

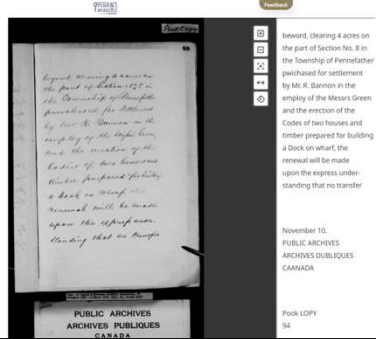
10

Initial Results: "Out of the Box"

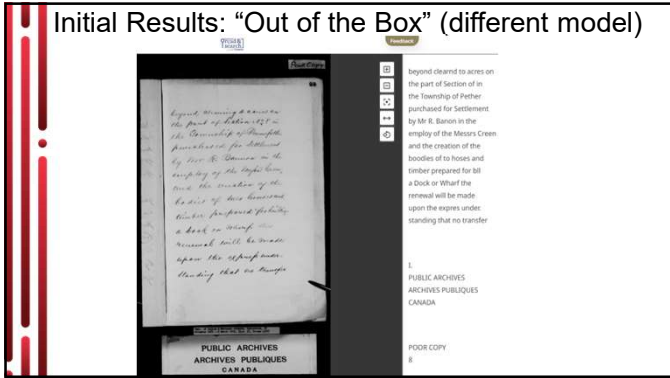


11

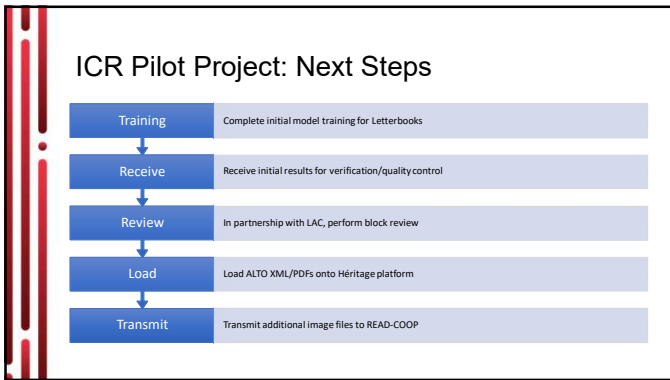
Initial Results: "Out of the Box"



12



13



14

How you can help

We still have approximately 1 million more images for the pilot, so suggestions are welcome!

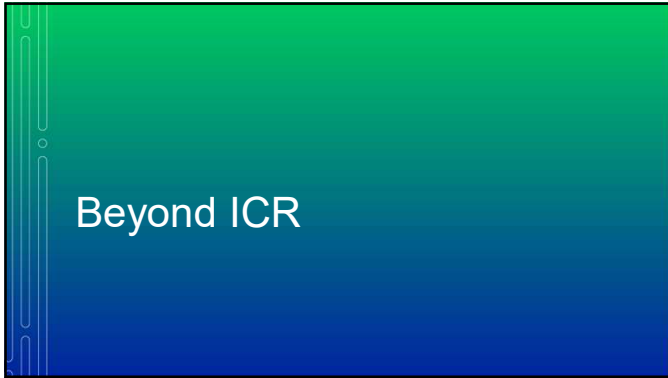
Identified series include:

- Annuity paylists
- Census records
- Black series
- Thousand series
- Red series

And if you have transcriptions, please share.

Email: jfriedman@crkn.ca

15



16

PDiiif: Custom PDF downloads from Héritage

- <https://pdiiif.jbaiter.de/>
- Instructions: <https://www.crkn-rcdr.ca/en/navigating-collections>
- Images only, no OCR

CRKN RCDR

17

Coming soon... improved search results display

Department of Indian Affairs : treaties, surrenders and agreements : T-9938

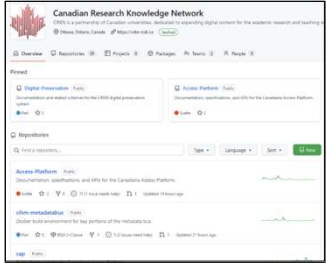
Previous result Next result 27/32 Images containing "treaty" in this reel

18

Future Enhancements

19

What is the Canadiana Platform today?



- Centralized web platform
- Monolithic slab of custom code
- Maintained by small team
- Technically open source code, available on GitHub

20

Replacing custom code with third-party solutions

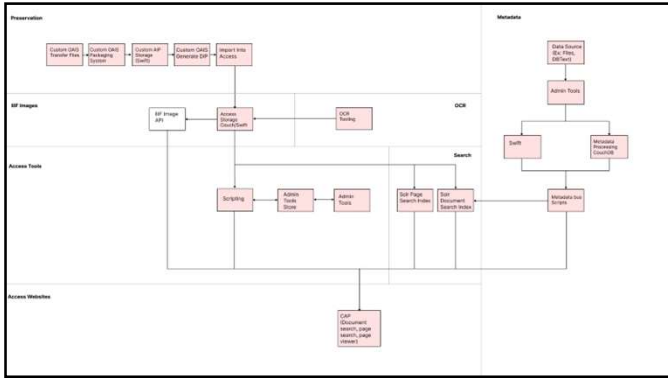
Archivematica for preservation

Blacklight for improved search and discovery

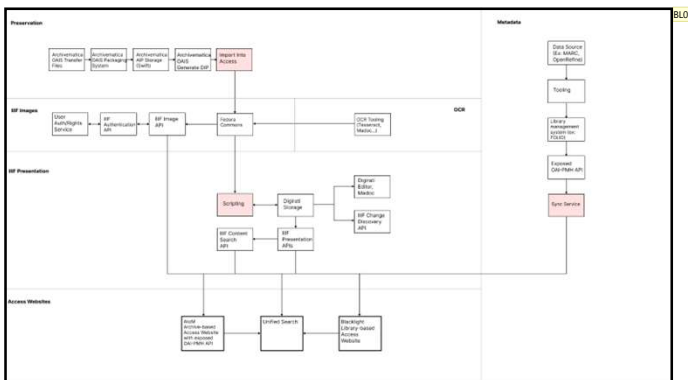
Digitati for collection management

FOLIO for better description

21



22



23

Canadian Research Knowledge Network
 Réseau canadien de documentation pour la recherche

Thank you! Questions?

jfriedman@crkn.ca

crkn-rodr.ca

24

Slide 23

- BLO** Russell said he would not be surprised if we could find a non-custom solution for the unified search at that point.
Brittney Lapierre, 2023-03-07T17:13:47.908