# Charity Scraper

CHARITY ORGANISATION SCRAPER

# USER MANUAL

Version 1.0

# Charity Scraper

Thank you for purchasing the Charity Organisation Scraper (further in text: COS). This software is meant for getting information out of Canadian Registry of Charity organisation. It has multi-process architecture and it uses Microsoft Access 2007 format database as its storing mechanism.

## System requirments.

Program is written using Microsoft .NET framework version 4.0, therefore minimum system requirments are as follows:

- Windows OS, at least xp sp3

- [.NET framework 4.0](#)

- [Microsoft Access Connector](#) (in case the Microsoft Access is not installed on the workstation)

## THE FILE STRUCTURE

Program needs multiple files in its folder in order for it to work properly.

**CharityScrape.exe.config –** holds the connection string to database.

**Database11_template.accdb** The database template that holds structure for the datastore database.

**chromedriver.exe** Chrome browser that is used to re-fetch failed records

**testMap.xml –** Used when fetching/updating the records. Holds directions to where each of the Organisation Detail Field is located in webpages. Also holds default url and directions for csv file download.

**testSelMap.xml –** Used when re-fetching the records. Holds directions to where each of the Organisation Detail Field is located in webpages. Also holds default url.

# OVERALL PROCESS.

When you first launch the application, it will create a file copy of database template and name it Database11.accdb. Further on this copy will be used as a storage mechanism for all the results obtained in the program work process.

Program has four modes of action:

## UPDATE DATA

Gathers information from the website, starting with the resultspage number specified. Then saves the results to database.

## SYNC WITH CSV

Downloads the csv information page and then updates record's real address information in database based on organisation registration number.

## RE-FETCH FAILED ROWS

Sometimes Update data might fail migth fail more than 3 times. In such times program skips the record and continues to next ones. That results in a record inside database that has no Organisation name. This action seeks such records and tries to obtain the information using alternative method – by fully rendering the page using chromedriver and then applying the update data process on it.
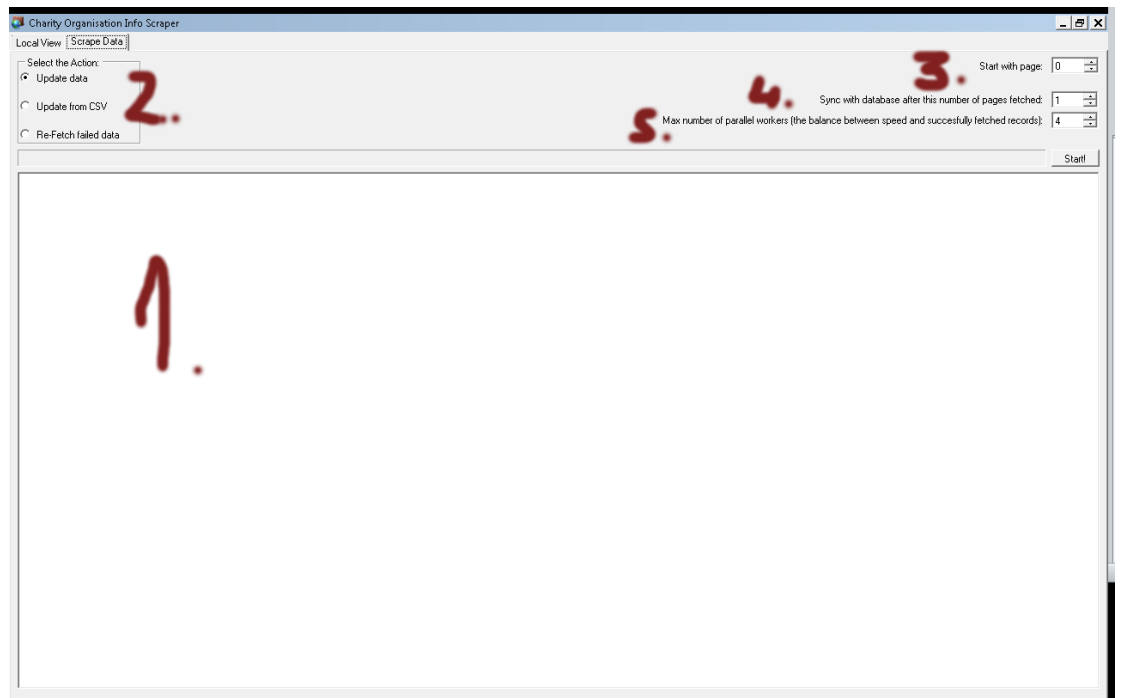
# The Local View tab

This section shows an overview of fetched results. It is filterable by organisation name. Clicking on the header of a column will sort the table by this column. Double clicking on any record will open the page from which this information was provided from. Clicking on an website column will open that particular link in a new system default browser window.

# The Scrape Data tab

This section allows to obtain the data from the Canadian Charity Organisation homepage and save it to database.



With (5) it is possible to specify maximum level of concurrent requests.

# XML MAPPINGS DESCRIPTION

The program needs two xml files for it to understand where the information is located in webpage. testMap.xml holds information for Update data process, while testSelMap holds same information for re-fetch failed rows process.

## LISTINGS PAGE MAPPINGS AND SOME SETTINGS

**FirstResultPageFormat** - Holds an URL to first page of advanced search results. Please notice that there must be a {0} written inside it in order for program to know where to put page number so to navigate between results pages.

**URLPrefix** - Holds Url that is prepended to every link in search results page so to open the corresponding details page. It must not include trailing slash (/)

NumHolder - Holds XPath to element that contains information about total records count in search results page.

**NumHolderRegEx** - Holds Regular Expression that is applied to NumHolder in order to extract number from its text.

ResultLinks - Holds XPath to elements (href attribute) that contains links to detail pages.

**DownloadResults** - Holds XPath to element that contains link (href attribute) to csv when geting addresses out of their csv.

## COMMON MAPPINGS

**Organization_name** - Holds XPath to element that contains information about the name of the organisation.

<!-- Match at groups[2] means quick info, groups[4] - Details page (like for 849923438RR0001)-->

**Organization_name_RegEx** - Holds Regular Expression that is applied to Organization_name in order to extract organisation name from its text.

## QUICK VIEW PAGE MAPPINGS

**Registration_no** - Holds XPath to element that contains information about the Registration number of the organisation.

**Designation** - Holds XPath to element that contains information about the Designation of the organisation.

**Web_site** - Holds XPath to element that contains information about the Website of the organisation.

**Programs_and_activities** - Holds XPath to element that contains information about the Programs and activities of the organisation.

**Status** - Holds XPath to element that contains information about the Status of the organisation.

**StatusDate** - Holds XPath to element that contains information about the StatusDate of the organisation.

**StatusDateFormat** - Holds XPath to element that contains information about the StatusDateFormat of the organisation.

**Operations_outside_Canada** - Holds XPath to element that contains number of Operations outside Canada of the organisation.

**Operations_outside_Canada_RegEx** - Holds Regular Expression that is applied to Operations_outside_Canada in order to extract number from its text.

**Total_revenue** - Holds XPath to element that contains information about the Total revenue of the organisation.

**Total_revenue_RegEx** - HHolds Regular Expression that is applied to Total revenue in order to extract number from its text.

**Total_expenses** - Holds XPath to element that contains information about the Total expenses of the organisation.

**Total_expenses_RegEx** - Holds Regular Expression that is applied to Total_expenses in order to extract number from its text.

**Total_compensation_for_all_positions** - Holds XPath to element that contains information about the Total_compensation_for_all_positions of the organisation.

**Full_time_employees -** Holds XPath to element that contains number of Full time employees in the organisation.

**Part_time_employees** - Holds XPath to element that contains number of  Part time employees in the organisation.

Charity Scraper

## DETAILS VIEW PAGE MAPPINGS

**DRegistration_no** - Holds XPath to element that contains information about the Registration number of the organisation.

**DDesignation** - Holds XPath to element that contains information about the Designation of the organisation.

**DWeb_site** - Holds XPath to element that contains information about the Web site of the organisation.

**DPrograms_and_activities** - Holds XPath to element that contains information about the Programs and activities of the organisation.

**DStatus** - Holds XPath to element that contains information about the Status of the organisation.

**DStatusDate** - Holds XPath to element that contains information about the StatusDate of the organisation.

**DStatusDateFormat** - Holds XPath to element that contains information about the StatusDateFormat of the organisation.

**DOperations_outside_Canada** - Holds XPath to element that contains number of Operations outside Canada of the organisation.

**DOperations_outside_Canada_RegEx** - Holds Regular Expression that is applied to DOperations_outside_Canada in order to extract number from its text.

**DTotal_revenue** - Holds XPath to element that contains information about the Total_revenue of the organisation.

**DTotal_revenue_RegEx** - Holds Regular Expression that is applied to Total_revenue in order to extract number from its text.

**DTotal_expenses** - Holds XPath to element that contains information about the Total_expenses of the organisation.

**DTotal_expenses_RegEx** - Holds Regular Expression that is applied to Total_expenses in order to extract number from its text.

**DTotal_compensation_for_all_positions** - Holds XPath to element that contains information about the Total compensation for all positions of the organisation.

**DFull_time_employees** - Holds XPath to element that contains information about the Full time employees of the organisation.

**DPart_time_employees** - Holds XPath to element that contains information about the Part time employees of the organisation.