# JournalList Ecosystem Web Crawler
Last Updated: July 4, 2022

This document describes how to use a Python based web crawler to download all the trust.txt files that can be reached by traversing all the trust.txt files referenced from a root URL (by default https://www.journallist.net/). This document is intended for Mac computers running macOS and assumes some familiarity with the "terminal" application and the UNIX shell.
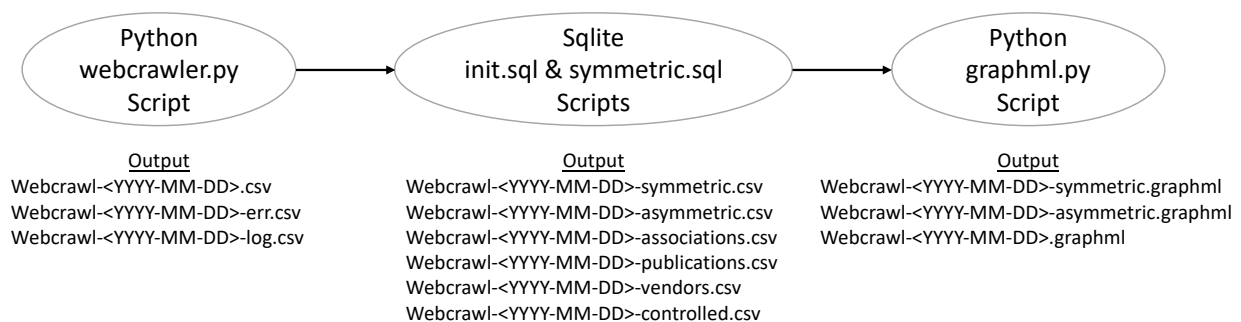
Preparing the environment:

1. Download the Python environment: https://www.python.org/downloads/
2. Install the Python requests module: https://docs.python-requests.org/en/master/user/install/#python-m-pip-install-requests
3. Download the yEd graphml editor/visualizer: https://www.yworks.com/products/yed/download#download

The cron.sh bash shell script captures this flow:
1. A Python webcrawler.py script crawls the trust.txt files generating .csv files capturing the results.
2. A pair of Sqlite scripts process the Webcrawl-<YYYY-MM-DD>.csv and Webcrawl-<YYYY-MM-DD>-err.csv files to derive the .csv files used in the final visualization step.
3. A Pytyon graphml.py script generates three .graphml files visualizing just the symmetric links found, just the asymmetric links found, and the complete social graph.

The processing flow is shown here.

| Python webcrawler.py Script | → | Sqlite init.sql & symmetric.sql Scripts | → | Python graphml.py Script |
|---|---|---|---|---|

| Output | Output | Output |
|---|---|---|
| Webcrawl-<YYYY-MM-DD>.csv | Webcrawl-<YYYY-MM-DD>-symmetric.csv | Webcrawl-<YYYY-MM-DD>-symmetric.graphml |
| Webcrawl-<YYYY-MM-DD>-err.csv | Webcrawl-<YYYY-MM-DD>-asymmetric.csv | Webcrawl-<YYYY-MM-DD>-asymmetric.graphml |
| Webcrawl-<YYYY-MM-DD>-log.csv | Webcrawl-<YYYY-MM-DD>-associations.csv | Webcrawl-<YYYY-MM-DD>.graphml |
|  | Webcrawl-<YYYY-MM-DD>-publications.csv |  |
|  | Webcrawl-<YYYY-MM-DD>-vendors.csv |  |
|  | Webcrawl-<YYYY-MM-DD>-controlled.csv |  |

In addition, there is a trust.txt QA python script that takes a trust.txt file as input and lists any errors encountered when parsing it.

| Input trust.txt | → | Python qa_trust_txt.py Script | → | Output List of errors whoislist.txt |
|---|---|---|---|---|

Running the webcrawler:

Open a terminal window and change directory to the Journalist directory. Pull the following files from the JournalList GitHub repository: https://github.com/brownwolf1355/JournalList

1. cron.sh - a bash shell script that runs the python webcrawler, processes the results through sqlite, and generates graphml files of the results.
2. webcrawler.py - a python script that recursively crawls trust.txt files to capture the state of the trust.txt ecosystem. It captures a copy of all of the trust.txt files it finds and generates a .csv file of the contents of all of them.
3. init.sql - the initialization sqlite script that creates the intermediate tables used in the following sql script.
4. symmetric - a sql script that generates .csv files containing the symmetric links in the trust.txt ecosystem and list of associations, publishers, and vendors discovered.
5. graphml.py - a python script that generates three graphml files containing the symmetric links, the asymmetric links, and the full ecosystem including both the symmetric and asymmetric links.
6. qa_trust_txt.py - a python script that parses a trust.txt file and lists any errors it contains.

Run the cron.sh shell script:

$ bash cron.sh

The shell and python scripts will create a directory "Webcrawl-YYYY-MM-DD" that will contain all the results of crawling the trust.txt files. This directory will contain the following:

1. The trust.txt files downloaded through the web crawling process. Each trust.txt file will be given a unique name "www.<domain>-trust.txt".
2. A log file useful for debugging, should anything go awry.
3. Several .csv files
   a. Webcrawl-YYYY-MM-DD.csv containing three columns: "srcurl," "attr", and "refurl". Representing the Source URL, Attribute, and Referenced URL on each line of all downloaded trust.txt files.
   b. Webcrawl-YYYY-MM-DD-err.csv containing four columns: "srcurl", "attr", "refurl", and "error". Representing the Source URL, Attribute, Referenced URL, and the Error encountered when the Referenced URL trust.txt file is attempted to be downloaded, a blank "www.<domain>-trust.txt" file will be created in these cases. Errors that are identified as "HTTP Status Code:" or "Content type: text/html" indicate that the server was reached, but a trust.txt file doesn't exist. Errors that are identified as "HTTP GET connection error exception occurred:" indicate that the server wasn't able to be reached and may indicate an error in the spelling of the Referenced URL.

c. Webcrawl-YYYY-MM-DD -symmetric.csv containing three columns: "srcurl," "attr", and "refurl". Representing the Source URL, Attribute, and Referenced URL for all symmetric links (e.g., "member" <-> "belongto" or "control" <-> "controlledby" or "vendor" <-> "customer"). The Attribute will be one of "member", "control", or "vendor".

d. Webcrawl-YYYY-MM-DD -asymmetric.csv containing three columns: "srcurl," "attr", and "refurl". Representing the Source URL, Attribute, and Referenced URL for all asymmetric links (e.g., "member" -> "belongto" reference without a corresponding "belongto" -> "member" reference). Note, these are the links listed in the Webcrawl-YYYY-MM-DD-err.csv file.

e. Webcrawl-YYYY-MM-DD-associations.csv containing one column "srcurl" listing all of the Source URLs that have more "member" attributes than "belongto" entries and all of the Referenced URLs that have "control" attributes.

f. Webcrawl-YYYY-MM-DD-publishers.csv containing one column "srcurl" listing all of the Source URLs that have more "belongto" attributes than "member" entries and all of the Referenced URLs that have "control" attributes.

g. Webcrawl-YYYY-MM-DD-vendors.csv containing one column "srcurl" listing all of the Referenced URLs that have "vendor" attributes.

h. Webcrawl-YYYY-MM-DD-symmetric.graphml containing the graphml diagram for the all symmetric links discovered.

i. Webcrawl-YYYY-MM-DD-asymmetric.graphml containing the graphml diagram for the all asymmetric links discovered.

j. Webcrawl-YYYY-MM-DD.graphml containing the graphml for the all links discovered, both symmetric and asymmetric.

k. Webcrawl-YYYY-MM-DD -links.json containing the JSON description of all the links for the ArangoDB graph analysis

l. Webcrawl-YYYY-MM-DD -urls.json containing the JSON description of all the urls for the ArangoDB graph analysis

Visualizing the results:

There are a few ways to analyze these results:

1. Through a yEd visual layout of the trust.txt ecosystem social graph
2. Through an Excel Pivot Table
3. Through a social graph analysis using ArangoDB

Visualizing the Social Graph with yEd:

To visualize the trust.txt ecosystem social graph, open the graphml files int yEd. Because the generated graphml files do not have any imbedded layout, the graph will be displayed with all the nodes on top of each other:

Use yEd's layout tool to display the trust.txt social graph. I have found that the "Organic" layout with "Preferred Edge Length" set to 100 to be most useful. You may wish to manually place the graph legend (consisting of 4 nodes and 4 edges):
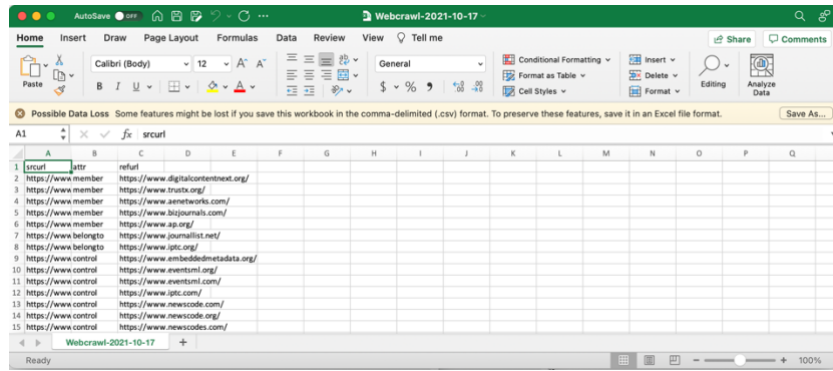


Use the selection tool to select these nodes and edges. You can obtain a better layout using the "Selection (Partial)" layout tool:
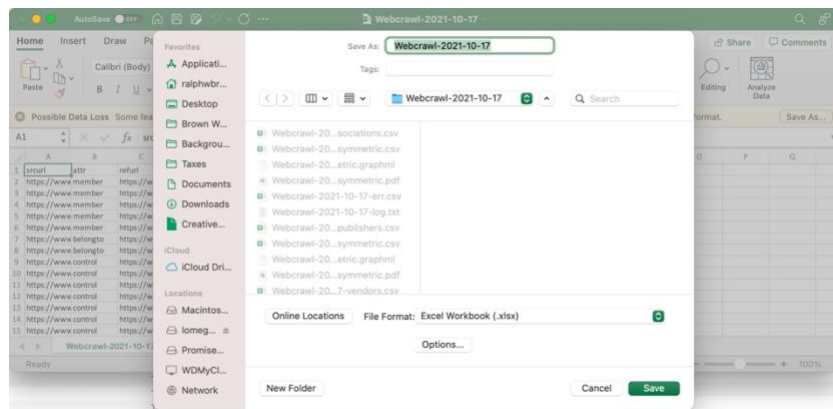


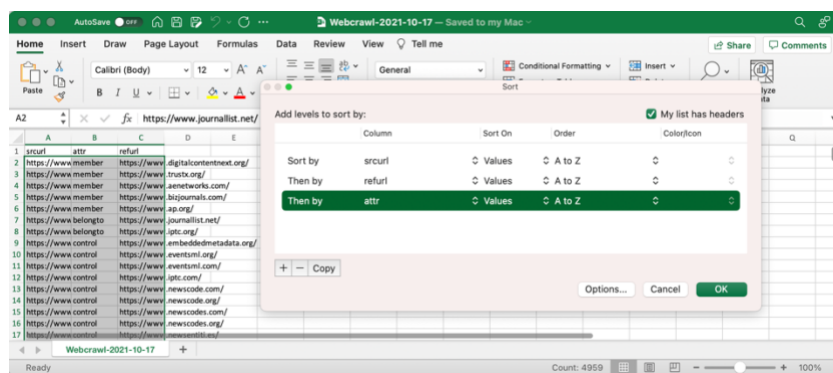The "Component" layout tool has options for organizing the layout of the social graph.

Analysis Using an Excel Pivot Table:

To create the Excel Pivot Table, open the file "Webcrawl-YYYY-MM-DD.csv" with Excel. It will notify you regarding Possible Data Loss.

Save the file as an Excel Workbook "Webcrawl-YYYY-MM-DD.xlsx".



Select all the contents of this worksheet and custom sort the contents, first by srcurl, then by refurl, and then by attr.



Select all again and insert a PivotTable using selected range and placing it on a new worksheet.

Add srcurl to the Rows of the PivotTable Fields and refurl below this. Add attr to the Columns of the PivotTable Fields and to the Sum Values.



Collapse all the sub-rows by selecting cell A5 and "Collapse Field" under "PivotTable Analyze". This will provide a summary of the different attributes used in each of the trust.txt files found.

Add a cell to capture the filename of the Excel workbook and create a title for the using the following formulas:

- Cell Q1
- Cell A1 =CONCAT("Status of the JournalList ecosystem ",MID(Q1,FIND("Webcrawl-",Q1)+9,10))

Create a list of the organization URLs with that have published a trust.txt file by placing the following in cell M5 of the piviot table:

=UNIQUE('[Webcrawl-YYYY-MM-DD.xlsx]Webcrawl- YYYY-MM-DD!A2:A<last row>,FALSE,FALSE)



Copy this cell to the two adjacent cells

Put headers above these in cells M4, N4, and O4 using the following:

| List of existing trust.txt urls | List of Attrs | List of referenced trust.txt urls |
|---|---|---|



Above this in cell M1, N1, and O1 create the following headers:

| JournalList Members | Existing trust.txt files | Referenced trust.txt files |
|---|---|---|



The number of JournalList members can be selected from the pivot table, number of existing published trust.txt files is obtained by =COUNTA(M5#), the number of referenced, but not published trust.txt files is obtained by =COUNTA(O5#). The list of attributes used across all the existing trust.txt files is show in Column N.

## Social Graph Analysis Using ArangoDB:

To create a social graph analysis using ArandoDB (https://www.arangodb.com/) first create a graph database for the webcrawl:

Login to system DB and create a new graph database for the webcrawl.



Login to the new DB and create a "links" edge collection and a "urls" document collection.

Upload the Webcrawl-YYYY-MM-DD -links.json and Webcrawl-YYYY-MM-DD -urls.json to their respective collections:



Then create a graph "urlrank" with Edge definitions from the "links" collection and the from/toCollections from the "urls" collection. If this is the first time, import the Webcrawl queries.



Run the queries "01 - Set link symmetry and weight" through "06 - Compute url ranking" in sequence. The url ranking can be generated by running the "List url rank by average" query and saving the CSV file.

The resulting CSV file can be imported into Excel and filtered as desired.

Generating trust.txt files from a website:

There are two options for running the scrapesite.py python script:

$ python3.10 sitescrape.py "Publication Name" **Error! Hyperlink reference not valid.**

This generates an output file <website>-trust.txt containing the trust.txt file for the referenced website.

$ python3.10 sitescrape.py input.csv output.csv

Where input.csv contains a list of publication names and their websites.
- input.csv - a .csv file containing a list of sites to scan [Name, Website]
- output.csv - the output .csv file containing [Website, Contact, Facebook, Instagram, Twitter, Youtube, LinkedIn, Vendor, Copyright, Control]

This output.csv can be run through an awk script to generate trust.txt files based on this output. The awk script will need to be modified for the specific Association. This example is for the Texas Press Associaiton.

$ awk -F "," -f tpa.awk tpa-output.csv

Open the output.csv file in Excel and ad the following columns M through R (where 339 is replaced by the number of rows in the sheet. Copy the equations through all the rows in the sheet and filter all of the columns.

| Start | End | Subdirectory | Copyright | Top URL | Count |
|-------|-----|--------------|-----------|---------|-------|
| =FIND("//" ,B2)+2 | =IF(RIGHT(B2,1)="/",LEN (B2)-1,LEN(B2)) | =IFERROR(FIND("/",MID(B 2,M2,N2-M2)),0) | =COUNTIF(K$2: K$339,K2) | =IF(O2<>0,LEFT(B2, M2+O2-1),"") | =IF(O2<>0,COUNTIF(Q$ 2:Q$339,Q2),0) |

Filter on column R ("Count") to select all values are not zero or one. This will generate a list of publications that published as subdirectories of a common website. You can then filter on column Q ("Top URL") for each website that has multiple publications. Check to see that there has been a trust.txt file generated for the Top URL. If not, you can use the scrapesite.py Python script to generate it as described above.

All the unique social media and contact sites for all these publications must be added to the trust.txt file of the Top URL. You can then delete the trust.txt files that were generated for the

subdirectories. They will be of the form "<website>-<subdirectory>-trust.txt". You should have one "<website>-trust.txt" file. Do this for each Top URL.

It may be necessary to generate a trust.txt file for publications with common ownership or control that are published on different websites, as opposed to through subdirectories of a common website. In this case you can filter on column L ("Control") and generate a trust.txt file for each controlling website. You then need to add "control=" entries for each website in column B ("Website").

<u>Running the QA tool:</u>

Run the qa_trust_txt.py python script:

$ python3.10 qa_trust_txt.py *filename*

The *filename* is a local trust.txt file to be checked.

For example, here is a sample trust.txt file:

# Brown Wolf Consulting LLC test trust.txt file

belongto=https://journallist.net/

controleby=https://www.brownwolfconsulting.com

member=https://www.brownwolfconsulting.org/

control=https://www.google.com/

social=https://www.linkedin.com/in/ralphwbrown/

social=https://twitter.com/RalphWBrown/

social=https://www.facebook.com/ralph.brown.5220/

The output looks like:

Invalid attribute at line 3 : controleby

Error at line 4 :  member=https://www.brownwolfconsulting.org/ - HTTP GET connection error

exception occurred: HTTPSConnectionPool(host='brownwolfconsulting.org', port=443): Max

retries exceeded with url: /trust.txt (Caused by

NewConnectionError('<urllib3.connection.HTTPSConnection object at 0x107d849a0>: Failed to

establish a new connection: [Errno 8] nodename nor servname provided, or not known'))

Error at line 5 :  control=https://www.google.com/ - HTTP Status Code: 404

It also generates a list of domain names in the whoislist.txt file that have connection errors. These can be used to check the status of these domains to understand the error further. For example:

$ whois brownwolfconsulting.org

Will generate the following:

% IANA WHOIS server

% for more information on IANA, visit http://www.iana.org

% This query returned 1 object

refer:        whois.pir.org

domain:      ORG

organisation: Public Interest Registry (PIR)

address:        11911 Freedom Drive 10th Floor,

address:        Suite 1000

address:        Reston, VA 20190

address:        United States

contact:       administrative

name:          Director of Operations, Compliance and Customer Support

organisation: Public Interest Registry (PIR)

address:        11911 Freedom Drive 10th Floor,

address:        Suite 1000

address:        Reston, VA 20190

address:        United States

phone:         +1 703 889 5778

fax-no:      +1 703 889 5779

e-mail:      ops@pir.org


contact:     technical

name:        Senior Director, DNS Infrastructure Group

organisation: Afilias

address:     Building 3, Suite 105

address:     300 Welsh Road

address:     Horsham, Pennsylvania 19044

address:     United States

phone:       +1 215.706.5700

fax-no:      +1 215.706.5701

e-mail:      tld-tech-poc@afilias.info


nserver:     A0.ORG.AFILIAS-NST.INFO 199.19.56.1 2001:500:e:0:0:0:0:1

nserver:     A2.ORG.AFILIAS-NST.INFO 199.249.112.1 2001:500:40:0:0:0:0:1

nserver:     B0.ORG.AFILIAS-NST.ORG 199.19.54.1 2001:500:c:0:0:0:0:1

nserver:     B2.ORG.AFILIAS-NST.ORG 199.249.120.1 2001:500:48:0:0:0:0:1

nserver:     C0.ORG.AFILIAS-NST.INFO 199.19.53.1 2001:500:b:0:0:0:0:1

nserver:     D0.ORG.AFILIAS-NST.ORG 199.19.57.1 2001:500:f:0:0:0:0:1

ds-rdata:    26974 8 2

4fede294c53f438a158c41d39489cd78a86beb0d8a0aeaff14745c0d16e1de32


whois:       whois.pir.org


status:      ACTIVE

remarks:     Registration information: http://www.pir.org


created:     1985-01-01

changed:     2020-10-27

source:     IANA


# whois.pir.org


NOT FOUND

>>> Last update of WHOIS database: 2021-10-21T15:10:55Z <<<