**JournalList Ecosystem Web Crawler**
Last Updated: October 17, 2021

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

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:
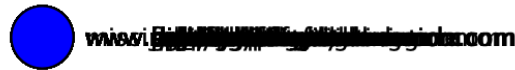
1. The trust.txt files downloaded through the web crawling process. Each trust.txt file will be given a unique name "www.<domain>-trust.txt".

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.
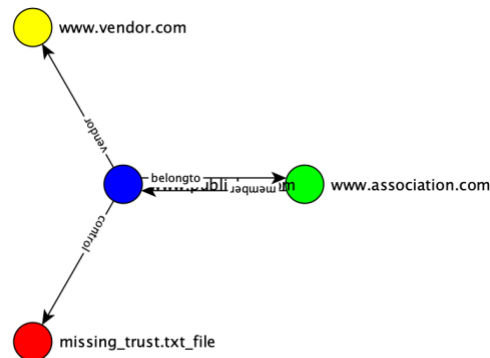
There are a couple ways to analyze these results:

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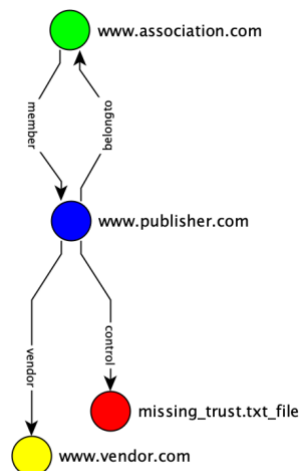2. Through an Excel Pivot Table.

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):
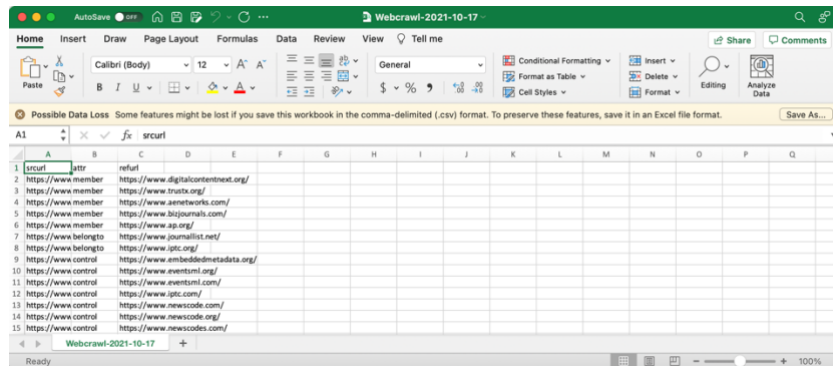


Use the selection tool to select these nodes and edges. You can obtain a better layout using the "Selection (Partial)" layout tool:



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

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.

**Status of the JournalList ecosystem 2021-10-17**

| JournalList Members | Existing trus | Referenced trust.txt files |
|---|---|---|
| 25 | 47 | 1397 |

| Count of attr | Column Labels | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Row Labels | belongto | contact | control | controlledby | member | social | vendor | Grand Total |
| https://www.aan.org/ | | 1 | 1 | | 87 | 3 | | 92 |
| https://www.anjournal.com/ | 1 | 1 | | | | | 2 | 4 |
| https://www.ap.org/ | 2 | 1 | 20 | | | 9 | | 32 |
| https://www.apnews.com/ | | 1 | | 1 | | | | 2 |
| https://www.archercountynews.com/ | 1 | 1 | | | | 2 | 2 | 6 |
| https://www.bctribune.com/ | 1 | 1 | | | | 1 | 2 | 5 |
| https://www.billtrack50.com/ | 2 | 1 | | | 4 | | | 7 |
| https://www.brownwolfconsulting.com/ | 1 | | | | 3 | | | 4 |
| https://www.canadianrecord.com/ | 1 | 1 | | | | 3 | 2 | 7 |
| https://www.colabnews.co/ | 3 | | | | 154 | 3 | | 160 |
| https://www.coloradofoic.org/ | 4 | 1 | 1 | | 24 | 4 | | 34 |
| https://www.coloradopressassociation.com/ | 7 | 1 | 1 | | 137 | 4 | | 150 |
| https://www.countystarnews.com/ | 1 | | | | | | 2 | 3 |
| https://www.dailycaller.com/ | 1 | 1 | 3 | | | 4 | | 9 |
| https://www.digitalcontentnext.org/ | | | | 1 | 65 | 2 | | 68 |
| https://www.doublemountainchronicle.com/ | 1 | 1 | | | | 2 | 2 | 6 |
| https://www.fastcompany.com/ | 3 | 1 | | 1 | 7 | | | 12 |
| https://www.fayettecountyrecord.com/ | 1 | 1 | 1 | | | 3 | 2 | 8 |
| https://www.firstdraftnews.org/ | 1 | 1 | | | 6 | | | 8 |
| https://www.flpress.com/ | 1 | 1 | 6 | | 113 | 2 | | 123 |
| https://www.giddingstimes.com/ | 1 | | | | | | 2 | 3 |
| https://www.haysfreepress.com/ | 4 | 1 | | | | 6 | | 11 |
| https://www.haysnewsdispatch.com/ | 4 | 1 | | | | 6 | | 11 |
| https://www.herefordbrand.com/ | 1 | | | | | | 2 | 3 |
| https://www.hernandosun.com/ | 4 | 1 | | | | 4 | | 9 |

| List of existing trust.txt urls | List of Attrs | List of referenced trust.txt urls |
|---|---|---|
| https://www.aan.org/ | contact | https://www.aan.org/contact-us... |
| https://www.anjournal.com/ | member | https://www.alibi.com/ |
| https://www.ap.org/ | social | https://www.austinchronicle.com/... |
| https://www.apnews.com/ | belongto | https://www.bendsource.com/ |
| https://www.archercountyne | vendor | https://www.bohemian.com/ |
| https://www.bctribune.com/ | control | https://www.boulderweekly.com/ |
| https://www.billtrack50.com/ | controlledby | https://www.c-ville.com/ |
| https://www.brownwolfconsulting.com/ | | https://www.charlestoncitypaper.r... |
| https://www.canadianrecord.com/ | | https://www.chicagoreader.com/ |
| https://www.colabnews.co/ | | https://www.chico.newsreview.co |
| https://www.coloradofoic.org/ | | https://www.chronogram.com/ |
| https://www.coloradopressassociation.co | | https://www.citybeat.com/ |
| https://www.countystarnews.com/ | | https://www.cityweekly.net/ |
| https://www.dailycaller.com/ | | https://www.clevescene.com/ |
| https://www.digitalcontentnext.org/ | | https://www.cltampa.com/ |
| https://www.doublemountainchronicle.co | | https://www.csindy.com/ |
| https://www.fastcompany.com/ | | https://www.cvindependent.com/ |
| https://www.fayettecountyrecord.com/ | | https://www.dailymemphian.com... |
| https://www.firstdraftnews.org/ | | https://www.dallasobserver.com/ |
| https://www.flpress.com/ | | https://www.digboston.com/ |
| https://www.giddingstimes.com/ | | https://www.eastbayexpress.com... |
| https://www.haysfreepress.com/ | | https://www.eriereader.com/ |
| https://www.haysnewsdispatch.com/ | | https://www.eugeneweekly.com/ |
| https://www.herefordbrand.com/ | | https://www.experiencecolumbias... |
| https://www.hernandosun.com/ | | https://www.facebook.com/altwe... |

Sheet1 | Webcrawl-2021-10-17