Social Feed Manager

Digital Preservation Coalition Web Archiving & Preservation Webinar Series February 6, 2020

Dan Kerchner kerchner@gwu.edu @DanKerchner Laura Wrubel lwrubel@gwu.edu @liblaura

Agenda

- Project background
- Current uses
- Collecting from Twitter
- SFM walkthrough
- Ethical and technical considerations

Social Feed Manager software

- Open source software by GW Libraries.
- User interface for collecting, managing, and exporting social media data.
- Collect from Twitter, Tumblr, Flickr, Sina Weibo.
- Libraries run this for their users as a service.
 (Not typically a local install on your laptop.)

More: <u>go.gwu.edu/sfm</u>

Project Background

- Initial development: 2013
- Version 1.0: 2016
- Grant funding:
 - IMLS
 - National Historical Publications & Records Commission (NARA)
 - CEAL







Institutions using SFM include

THE GEORGE WASHINGTON UNIVERSITY

WASHINGTON, DC

Stanford University







RADCLIFFE INSTITUTE FOR ADVANCED STUDY HARVARD UNIVERSITY

Social media research



INFORMATION, COMMUNICATION & SOCIETY, 2017

Populist communication by digital means: presidential Twitter

https://doi.org/10.1080/1369118X.2017.1328521

VOL. 20, NO. 9, 1330-1346

Routledge

Taylor & Francis Group

(R) Check for updat

GW Arts & Sciences

Following

 \checkmark

Congratulations to CCAS alums Sally Nuamah, BA '11, and Tara Dorfman, BA '11, who were both named to Forbes 2019 30 Under 30 lists! #GWCCAS #GWU



11:29 AM - 2 Dec 2018



```
created_at: "Sun Dec 02 16:29:00 +0000 2018",
                                                                                            in reply to status id: null,
  id: 1069267199318216700,
                                                                                            in reply to status id str: null,
  id str: "1069267199318216704",
                                                                                            in reply to user id: null,
  full text: "Congratulations to CCAS alums Sally Nuamah, BA '11, and Tara Dorfman
                                                                                           in reply to user id str: null,
  named to Forbes 2019 30 Under 30 lists! #GWCCAS #GWU https://t.co/kPWyKrJ9Ux",
                                                                                            in reply to screen name: null,
                                                                                          - user: {
  truncated: false,
                                                                                               id: 54639174,
+ display text range: [...],
                                                                                               id str: "54639174",
- entities: {
                                                                                               name: "GW Arts & Sciences",
    - hashtags: [
                                                                                               screen name: "gwucolumbian",
       - {
                                                                                               location: "Washington, D.C.",
              text: "GWCCAS",
                                                                                               description: "Official Twitter account for the Columbian College of Arts and
           - indices: [
                                                                                               Sciences at The George Washington University. Home of the Engaged Liberal
                 132,
                                                                                               Arts.",
                 139
                                                                                               url: "https://t.co/g0ys0T59mJ",
                                                                                             - entities: {
                                                                                                 - url: {
          },
                                                                                                    - urls: [
        - {
                                                                                                        - {
              text: "GWU",
                                                                                                             url: "https://t.co/q0ys0T59mJ",
            - indices:
                                                                                                             expanded url: "http://columbian.gwu.edu",
                 140,
                                                                                                             display_url: "columbian.gwu.edu",
                 144
                                                                                                           - indices: [
                                                                                                                 0,
              1
                                                                                                                 23
      ],
      symbols: [ ],
      user mentions: [ ],
                                                                                                   },
      urls: [ ],
                                                                                                 - description: {
    - media: [
                                                                                                      urls: []
       - {
                                                                                                   }
                                                                                               },
              id: 1069267196357091300,
                                                                                               protected: false,
              id str: "1069267196357091328",
                                                                                               followers count: 4322,
           - indices: [
                                                                                               friends count: 608,
                 145,
                                                                                               listed count: 136,
                 168
                                                                                               created_at: "Tue Jul 07 18:49:10 +0000 2009",
              1,
                                                                                               favourites count: 4202,
             media url: "http://pbs.twimg.com/media/DtbM2ZDXoAAMvQ7.jpg",
                                                                                               utc offset: null,
             media url https: "https://pbs.twimg.com/media/DtbM2ZDXoAAMvQ7.jpg",
                                                                                               time zone: null,
              url: "https://t.co/kPWvKrJ9Ux",
                                                                                               qeo enabled: true,
             display url: "pic.twitter.com/kPWyKrJ9Ux",
                                                                                               verified: false,
                                                                                               statuses count: 9273,
              expanded url: "https://twitter.com/gwucolumbian/status/10692671993182
                                                                                               lang: "en",
              type: "photo".
```

Collecting from Twitter's APIs

- SFM uses the free standard APIs
- API responses are JSON
- SFM handles rate limiting, authentication, organizing collections
- Requires one set of Twitter developer keys for the application and Twitter authentication for each user collecting data.

Twitter collection types

- User timeline: up to 3,200 tweets in the past, for each account
- Search: sample of tweets from the past 7-10 days
- Filter stream: real-time, forward-looking only
- Sample stream ~ 0.5% sample of all tweets; ~2GB per day

Collecting from Twitter: Data Extracts

- JSON
- CSV and Excel (selected fields from JSON)
- Filtered by date range, screen name, etc.

CSV/Excel extract

en

FALSE

1

id	tweet_url	created_a	t parsed_c	re user_sc	ree text						tweet_type	coordin	hashtags	media	urls	favorite_cc i
106970559017435	https://twitte	e Mon Dec (2018-12-0)3 gwucolu	umt Kyral	h Altman, a senio	r majoring in hu	man services a	nd social jus	tice, was a	original		CCASOutFr	https://tv	vitter.com/gv	v 0
106967441805249	https://twitt	e Mon Dec (2018-12-0)3 gwucolu	umł RT @	CorcoranGW: Pr	ofessor Maria de	l Carmen Mon	toya, featur	ed in a rec	retweet					0
106934420665575	https://twitte	e Sun Dec 02	2 2018-12-0)2 gwucolu	umt The r	new #GWCCAS In	ternational Budo	ly Program wai	nts to help i	nternation	original		GWCCAS G	i https://tv	vi https://bit	. 1
106926719931821	https://twitte	e Sun Dec 02	2 2018-12-0)2 gwucolu	umt Cong	gratulations to CC	AS alums Sally N	uamah, BA '11	, and Tara D	orfman, B/	original		GWCCAS G	i https://tv	vitter.com/gv	/ 9
106896923730842	https://twitte	e Sat Dec 01	2018-12-0)1 gwucolu	umł A #G	WCCAS junior an	d her peers got t	he opportunity	to see Mic	helle Oban	original		GWCCAS C	https://tv	vi https://bit	. 3
										1.0				1 11.		
106888619715248	https://twitte	e Sat Dec 01	2018-12-0)] gwucoli	umt Last v	week, @Corcora	GW hosted Was	shington D.C.'s	annual Nati	onal Portfo	original		GWU GWC	https://tv	vi https://bit	. 3
in_reply_tc in_reply			. 2018-12-0 place			week, @Corcorai et_ccretweet_or				user_id						. 3 lov user_frien
					_seretwee	· -		quc retweet_o		user_id	user_creat	user_de	fau user_de	escr user_fa		lov user_frien
		lang		possibly_	_seretwee	et_ccretweet_or_ 0		quc retweet_o	source	user_id 54639174	user_creat Tue Jul 07	user_de	fat user_de Official	escruser_fa	avou user_fol	lov user_frien 22 608
		lang en		possibly_	_s(retwee	et_ccretweet_or_ 0	qu(retweet_or_	quc retweet_o	source <a h1<br="" href="ht</td><td>user_id
54639174
54639174</td><td>user_creat
Tue Jul 07
Tue Jul 07</td><td>user_de
FALSE
FALSE</td><td>fat user_de
Official
Official</td><td>escr user_fa
Twi 4
Twi 4</td><td>avou user_fol
202 43</td><td>lov user_frien
22 608
22 608</td></tr><tr><td></td><td></td><td>en
en</td><td></td><td>possibly_
FALSE</td><td>_s(retwee</td><td>et_ccretweet_or_
0</td><td>qu(retweet_or_</td><td>quc retweet_o</td><td>source
							

<a href="h154639174 Tue Jul 07 : FALSE Official Twi

608

4202

4322

SFM Walkthrough

Social Feed Manager empowers researchers and archivists to build collections of social media data from multiple platforms.

Log In

If you have not created an account yet, then please sign up first.

Username	Flickr, and Sina Weibo. Read more about Social Feed Manager
Username	here.
Password	
Password	
Forgot your password?	
🗌 Remember me	
Log in	

Collecting and using data from social media platforms is subject to those platforms' terms (Twitter, Flickr, Sina Weibo, Tumblr), as you agreed to them when you created your social media account. Social Feed Manager respects those platforms' terms as an application (Twitter, Flickr, Sina Weibo, Tumblr).

Social Feed Manager empowers researchers and archivists to create collections of social media data from Twitter, Tumblr,

Social Feed Manager provides data to you for your research and academic use. Social media platforms' terms of service **generally do not allow republishing of full datasets**, and you should refer to their terms to understand what you may share. Authors typically retain rights and ownership to their content.

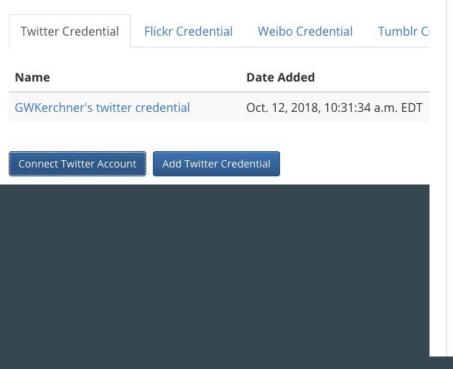
In addition to respecting the platforms' terms, as a user of Social Feed Manager and data collected within it, it is your responsibility to consider the ethical aspects of collecting and using social media data. Your discipline or professional organization may offer guidance. Here are some ethical and privacy guidelines you may want to consider.

tials Exports

Welco

Credentials

Credentials are used to authorize Social Feed Manager to collect data from T Authorize Social Feed Manager by connecting your account or adding creden



y

Monitor

Authorize Social Feed Manager (GW Sandbox) to use your account?

Username or email	
Password	
Remember me · Forgot	nase



This application will be able to:

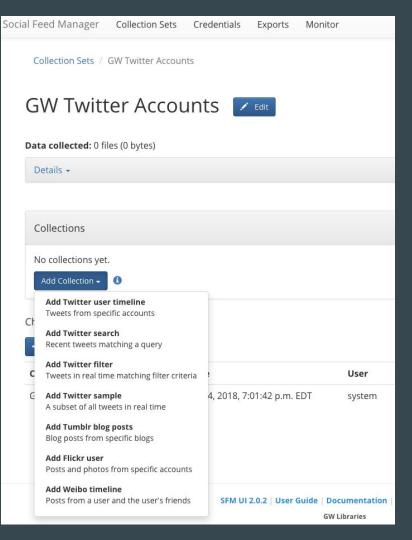
- · Read Tweets from your timeline.
- · See who you follow.

Will not be able to:

- · Follow new people.
- · Update your profile.
- · Post Tweets for you.
- Access your direct messages.
- · See your email address.
- · See your Twitter password.

Sign up for Twitter >





Collection Sets / GW Twitter Accounts / Add Twitter user timeline

Add Twitter user timeline

* indicates a required field.

Collection name*

GWU official accounts

Description

Public link

Link to a public version of this collection, e.g., in a data repository.

Credential*

GWKerchner's twitter credential

Incremental harvest

Only collect new items since the last data retrieval.

Automatically delete seeds for deleted / not found accounts.

□ Automatically delete seeds for suspended accounts.

Automatically delete seeds for protected accounts.

Schedule*

Every week

How frequently you want data to be retrieved.

End date

× III

If blank, will continue until stopped.

Sharing*

Group only

Who else can view and export from this collection. Select "All other users" to share with all Social Feed Manager users.

Change Note

Further information about this addition.

Save Cancel

Collection Sets / GW Twitter Accounts / GWU official accounts / Request Export

Request Export

Seed choice*

All seeds

 \bigcirc Active seeds only

 \bigcirc Selected seeds only

None selected +

Export format*

Excel (XLSX)	•
Maximum number of items per file	

Maximum number of items per file

	.000

Deduplicate (remove duplicate posts)

imit by item date range		
ltem date start		
	×	
ltem date end		
	×	

The timezone for dates entered here are America/New_York. Adjustments will be made to match the time zone of the items. For example, dates in tweets are UTC.

Limit by harvest date range

Harvest date start

		×	
Harves	st date end		
		×	

Collecting and using data from social media platforms is subject to those platforms' terms (Twitter, Flickr, Sina Weibo, Tumblr), as you agreed to them when you created your social media account. Social Feed Manager respects those platforms' terms as an application (Twitter, Flickr, Sina Weibo, Tumblr).

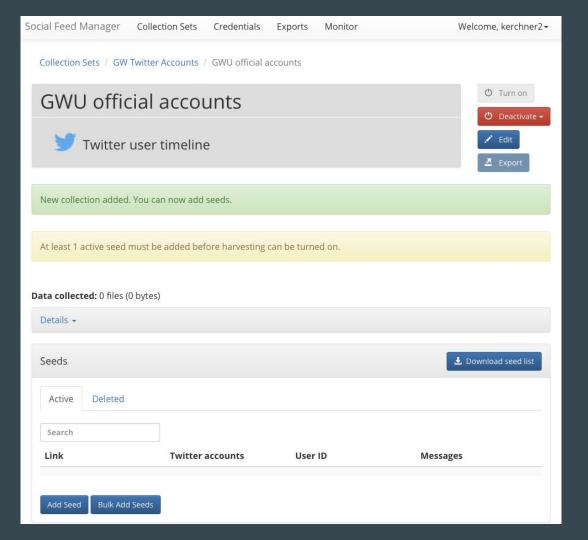
Social Feed Manager provides data to you for your research and academic use. Social media platforms' terms of service **generally do not allow republishing of full datasets**, and you should refer to their terms to understand what you may share. Authors typically retain rights and ownership to their content.

In addition to respecting the platforms' terms, as a user of Social Feed Manager and data collected within it, it is your responsibility to consider the ethical aspects of collecting and using social media data. Your discipline or professional organization may offer guidance. Here are some ethical and privacy guidelines you may want to consider.

SFM UI 2.0.2 | User Guide | Documentation | Citing | Contact Us

× III

GW Libraries



Collection Sets / GW Twitter Accounts / GWU official accounts / Add Twitter user timeline seeds

Add Twitter user timeline seeds

Seeds type*

Screen Name

○ User id

Bulk Seeds*

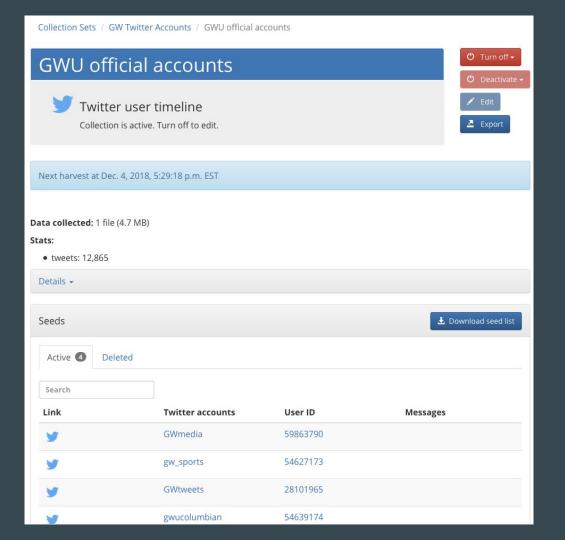
gw_sports GWmedia GWtweets

gwucolumbian

Enter each seed on a separate line.

Change Note

Further information about this addition.



Exports Monitor	Welcome, kerchner2 -
counts / Export	
Size	
2.0 KB	
2.8 MB	
ls in the export.	
	counts / Export Size 2.0 KB

Collecting and using data from social media platforms is subject to those platforms' terms (Twitter, Flickr, Sina Weibo, Tumblr), as you agreed to them when you created your social media account. Social Feed Manager respects those platforms' terms as an application (Twitter, Flickr, Sina Weibo, Tumblr).

Social Feed Manager provides data to you for your research and academic use. Social media platforms' terms of service **generally do not allow republishing of full datasets**, and you should refer to their terms to understand what you may share. Authors typically retain rights and ownership to their content.

In addition to respecting the platforms' terms, as a user of Social Feed Manager and data collected within it, it is your responsibility to consider the ethical aspects of collecting and using social media data. Your discipline or professional organization may offer guidance. Here are some ethical and privacy guidelines you may want to consider.

Collection Sets

A collection set is a group of collections around a particular topic or theme. Collections sets are active when there is at least one active collection within them. Collection sets are inactive when all collections have been deactivated and are no longer harvesting.

Active 22	Inactive 13	Shared 23	Other Active 109	Other Inactive 12	
Name		Collecti	ons	Date Added	Groups
115th U.S. Cor	ngress	3 collect	ions	Jan. 27, 2017, 10:47:	7:12 a.m. EST GW Libraries Scholarly Technology Group
2017-2020 Fee	deral Term	2 collect	ions	Jan. 20, 2017, 11:26:	6:54 a.m. EST GW Libraries Scholarly Technology Group
Alaska Earthquake		2 collections		Nov. 30, 2018, 1:37:.	7:25 p.m. EST GW Libraries Scholarly Technology Group
China Anti-Co	rruption	48 colle	ctions	June 29, 2016, 11:49	9:14 p.m. EDT CEAL Grant
Climate chang	je	1 collect	ion	Sept. 21, 2017, 8:11:	1:54 a.m. EDT GW Libraries Scholarly Technology Group
Corcoran		1 collect	ion	March 29, 2017, 5:1	17:54 p.m. EDT GW Libraries Scholarly Technology Group
Foreign leaders		1 collect	ion	March 11, 2018, 9:4	43:54 p.m. EDT GW Libraries Scholarly Technology Group
Governors		1 collect	ion	March 11, 2018, 9:34	34:03 p.m. EDT GW Libraries Scholarly Technology Group

Technical info

Requires:

- Twitter application credentials
- Linux server
- Storage

Social Feed Manager runs as a set of Docker containers

- Application written in Python
- Small database (postgresql) for user info and metadata
- Collected content is stored on the file server as WARC files

Ethical and privacy concerns

Social media data comes from people

- Consider the impact of your work on the creator of the social media. Be thoughtful collecting social media of vulnerable individuals.
- When publishing, get permission from creator for quotes from social media if possible, do not rely on anonymizing posts.
- Respect a user's deleting their content.
- Be familiar with platform terms of use. Don't republish full datasets and share in accordance with terms (e.g., tweet ids only)

SFM software

- <u>Open source on GitHub</u>
- Issues tracked in the sfm-ui repository: we welcome requests!
- Documentation for installation and use: <u>sfm.readthedocs.io/</u>
- Main project site (blog posts and more): <u>go.gwu.edu/sfm</u>

Thank you

Social Feed Manager sfm@gwu.edu <u>go.gwu.edu/sfm</u>

Dan Kerchner kerchner@gwu.edu @DanKerchner Laura Wrubel lwrubel@gwu.edu @liblaura

TweetSets tweetsets.library.gwu.edu

Select source datasets

Select the source datasets from which to create a new dataset.

	>
Hurricane Florence	
Hurricane Michael	
Solar Eclipse	
te User Timelines (22,251 tweets)	
ion 8,340,668 tweets	
er Timelines 11,846 tweets	
,626,254 tweets	
7,140 tweets	
Date Filter 3,182,243 tweets	
te User Timelines 52,395 tweets	
	Hurricane Michael Solar Eclipse te User Timelines (2,251 tweets) on (8,340,668 tweets) er Timelines (1,846 tweets) 626,254 tweets filde tweets pate Filter (3,182,243 tweets)