{"created_at": "Thu Apr 06 15:24:15 +0000 2017", "id_str": "850006245121695744", "text": "1\/ Today we\u201
9re sharing our vision for the future of the Twitter API platform!\nhttps:\/\/t.co\/XweGngmxlP", "user": {"
id": 2244994945, "name": "Twitter Dev", "screen_name": "TwitterDev", "location": "Internet", "url": "https:
\/\/dev.twitter.com\/", "description": "Your official source for Twitter Platform news, updates & events.
Need technical help? Visit https:\/\/twittercommunity.com\/ \u2328\ufe0f #TapIntoTwitter"}, "place": { },
"entities": { "hashtags": [ ], "urls": [ { "url": "https:\/\/t.co\/XweGngmxlP", "unwound": { "url": "https:
\/\/cards.twitter.com\/cards\/18ce53wgo4h\/3xo1c", "title": "Building the Future of the Twitter API
Platform"} } ], "user_mentions": [ ] } } {"created_at": "Thu Apr 06 15:24:15 +0000 2017", "id_str": "850006
245121695744", "text": "1\/ Today we\u2019re sharing our vision for the future of the Twitter API platform!
\nhttps:\/\/t.co\/XweGngmxlP", "user": {"id": 2244994945, "name": "Twitter Dev", "screen_name": "TwitterDev"
, "location": "Internet", "url": "https:\/\/dev.twitter.com\/", "description": "Your official source for
Twitter Platform news, updates & events. Need technical help? Visit https:\/\/twittercommunity.com\/ \u2328
\ufe0f #TapIntoTwitter"}, "place": { }, "entities": { "hashtags": [ ], "urls": [ { "url": "https:\/\/t.co\
/XweGngmxlP", "unwound": { "url": "https:\/\/cards.twitter.com\/cards\/18ce53wgo4h\/3xo1c", "title":
"Building the Future of the Twitter API Platform"} } ], "user_mentions": [ ] } } {"created_at": "Thu Apr 06
15:24:15 +0000 2017", "id_str": "850006245121695744", "text": "1\/ Today we\u2019re sharing our vision for
the future of the Twitter API platform!\nhttps:\/\/t.co\/XweGngmxlP", "user": {"id": 2244994945, "name":
"Twitter Dev", "screen_name": "TwitterDev", "location": "Internet", "url": "https:\/\/dev.twitter.com\/",

# {COLLECTION ETHICS}

# Collection Ethics

***Ethically curating collections of Twitter data***

When curating a collection of Twitter data, researchers have a legal responsibility to protect the privacy of the users represented in the dataset. Beyond legal limits placed on research, there is a [movement](#) towards more ethical collection and archival practices. [A 2018 study](#) found that over 60 percent of respondents were unaware that their tweets were sometimes used by researchers. This tutorial examines ethical challenges posed by the collection design, collection, and documentation processes explored in previous tutorials. [The STEP framework](#) is used as the ethical framework for curating Twitter data collections. While STEP provides a general framework, any researcher working with sensitive Twitter data should point all ethics-related questions to their Institutional Review Board (IRB). To contact UNLV's IRB, email: [IRB@unlv.edu](mailto:IRB@unlv.edu).

**Difficulty level:** Beginner

**Prerequisite(s)**
- [Tweet JSON](#)
- [Command Line](#)
- [Collection Design](#)
- [Collection with Twarc](#)
- [Collection Documentation](#)

**Lesson objectives**
- Review the STEP Framework for curating social media data
- Learn how to bowdlerize Tweets

**Key Terms**
- Twarc
  - A command line tool and python library for collecting Twitter JSON data
- Bowdlerization

- A form of purging anything deemed harmful or offensive from an artistic work, or other type of writing of media.

# Table of Contents

# Introduction to the STEP Framework

'STEP' was developed by Sara Mannheimer and Elizabeth Hull to provide researchers with an ethical framework for curating social media data. The framework is grounded in three guiding principles: Value analysis, responsibility, and continual inquiry. Depending on the collection subject, the tweets may contain information that presents a risk to the original creators of the content. As you move through this tutorial, keep in mind the ethical challenges presented by your own Twitter collection(s). The responsibility is on you, the researcher, to protect the privacy and safety of the users in your dataset. Review the framework below:

**Step framework**

**S -** Sensitive content
**T -** Transparent collection methods
**E -** Expectation of privacy
**P -** Platform policies

{"created_at": "Thu Apr 06 15:24:15 +0000 2017", "id_str": "850006245121695744", "text": "1\/ Today we\u2019re sharing our vision for the future of the Twitter API platform!\nhttps:\/\/t.co\/XweGngmxlP", "user": {"id": 2244994945, "name": "Twitter Dev", "screen_name": "TwitterDev", "location": "Internet", "url": "https:\/\/dev.twitter.com\/", "description": "Your official source for Twitter Platform news, updates & events. Need technical help? Visit https:\/\/twittercommunity.com\/ \u2328\ufe0f #TapIntoTwitter"}, "place": { }, "entities": { "hashtags": [ ], "urls": [ { "url": "https:\/\/t.co\/XweGngmxlP", "unwound": { "url": "https:\/\/cards.twitter.com\/cards\/18ce53wgo4h\/3xo1c", "title": "Building the Future of the Twitter API Platform"} } ], "user_mentions": [ ] } } {"created_at": "Thu Apr 06 15:24:15 +0000 2017", "id_str": "850006245121695744", "text": "1\/ Today we\u2019re sharing our vision for the future of the Twitter API platform!\nhttps:\/\/t.co\/XweGngmxlP", "user": {"id": 2244994945, "name": "Twitter Dev", "screen_name": "TwitterDev", "location": "Internet", "url": "https:\/\/dev.twitter.com\/", "description": "Your official source for Twitter Platform news, updates & events. Need technical help? Visit https:\/\/twittercommunity.com\/ \u2328\ufe0f #TapIntoTwitter"}, "place": { }, "entities": { "hashtags": [ ], "urls": [ { "url": "https:\/\/t.co\/XweGngmxlP", "unwound": { "url": "https:\/\/cards.twitter.com\/cards\/18ce53wgo4h\/3xo1c", "title": "Building the Future of the Twitter API Platform"} } ], "user_mentions": [ ] } }

# S

Is the information being studied of a **sensitive** nature?
Are the research subjects from vulnerable populations?

# T

Is there sufficient documentation to make the data resuseable and collection methods **transparent**?

# E

Did subjects have an **expectation** of privacy?
Was consent obtained for research and/or sharing?
Are the data properly anonymized, or can they be made so?

# P

Are the data keeping with the policies of the social media **platform**?

## Can the social media data be shared openly in a manner that is both safe and useful?

**Framework developed by Sara Mannheimer and Elizabeth Hull**

STEP

# Sensitive Content

While you may not be collecting around a particularly sensitive subject, it is still possible that information contained in the dataset could put the lives of the content creators in danger. Without an effective way to gather mass consent for research purposes, Twitter data must be treated with care. Ask yourself the following when in the collection design stage of the curation process:

1. Is the information being studied of a sensitive nature?
2. Are the research subjects from vulnerable populations?

If you answered yes to one or both of those questions, it is your responsibility as a researcher to take the necessary steps to protect users within your dataset. This includes:

- Sharing only the Tweet ids for your collection when sharing your collection outside of your institution
- Bowdlerizing any sensitive content published within a study.
    - Example: Removing all user handles and editing tweets so that, while the meaning is not lost, the tweets will not be returned in a search
- If collecting from a select group of users, you may want to consider requesting written consent from each

Bowdlerizing your tweets is a step you can take to protect the identity of the users in your collection. Below you will find an example of a tweet before and after it has been bowdlerized.

### Bowdlerizing Tweets

Mary and Alison Janet are attending a protest in their town. You are collecting data around the protest, and planning on publishing the results immediately for researchers on campus to access and analyze the data. View the tweet sent by Mary below:

**@janetmary:** **The new meeting point for the #Rally is 104 E. Marks St. My cell is 702-806-1817, txt with any questions. You can also reach out to @alisonjanet w/ questions.**


**Image attached to tweet**

There are several pieces of information in this tweet that should be modified to protect the identity of the user. Let's break them down.

## User handle
User handles are the most obvious identifier of a tweet. Unless you have a very specific use case, usually requiring user consent, publishing user handles in your work can present major ethical concerns.

## Addresses
All information that could be used to identify a person, like addresses, should be removed.

## Phone numbers
All information that could be used to identify a person, like phone numbers, should be removed.

## Images
If the subject of the photo is identifiable, the image should not be published without the consent of the user.

*Example of a Bowdlerized tweet*

**@user1:** **We have the new meeting point for the #Rally. Txt with any questions. You can also reach out to @user2 w/ questions.**

### User handle

The original sender of the tweet has a username that has been anonymized to protect their identity. The second user handle (@alisonjanet) was removed and replaced with a false name. The changed name does not alter the meaning of the original tweet.

### Addresses

The address was removed to protect the original poster of the tweet. This change did not alter the meaning of the original tweet.

### Phone numbers

The phone number was removed. This change did not alter the meaning of the original tweet.

### Images

The image was removed because it features two identifiable faces.

# Transparent Collection Methods

After creating a README file for a collection, ask yourself the following question:

1. Is there sufficient documentation to make the data reusable and collection methods transparent?

In the Collection Documentation tutorial, you completed a checklist and created a README file that will help you to make your collection process transparent to other researchers. Going forward, keep in mind the STEP framework when creating and updating your documentation. Consider whether or not the users in your dataset would face consequences if a bad actor were to gain access to your collection.

# Expectation of Privacy

The 2018 study that was referenced at the beginning of this tutorial stated that 42.7% of their respondents believed that researchers were not able to use their tweets without their permission. While this gap in knowledge may fall on the platform creators, the onus is on those who work with social media data to respect the user data they handle. Ask yourself the following three questions when collecting Twitter data:

1. Did subjects have an expectation of privacy?

2. Was consent obtained for research and/or sharing?
3. Are the data properly anonymized, or can they be made so?

Chances are high that you answered no to one or more of these questions. Due to the often macro scale of Twitter data, obtaining consent and confirming users' expectation of privacy can prove difficult. Just because you answered no to any of the three questions does not mean you need to dispose of your data—use these questions to motivate your work to protect the privacy of the users in your collection.

# Platform Policies

In the [Collection Design tutorial](#), you reviewed the Terms of Service and Developer Agreement for Twitter. Ask yourself the following question during the collection process:

1. Are the data keeping with the policies of the social media platform?

If you completed the Collection Design, Collection with Twarc, and Collection Documentation tutorials, your answer to that question should be an easy yes. By dehydrating your dataset to create a list of unique Twitter ids, you have taken the first step in the ethical sharing of your data. Including important disclaimers about the sensitive nature of the data in your README is another step towards protecting Twitter users' right to privacy. Finally, in this tutorial, the concept of bowdlerizing content was introduced. Editing any published tweets so that the original poster is unidentifiable will help to keep with the policies Twitter laid out in their Developer Agreement.