
Reliable and Efficient Long-Term Social Media Monitoring

Jian Cao

Postdoctoral Scholar In Election Integrity and Data Science
Division of Humanities and Social Sciences
California Institute of Technology
Pasadena, CA 91105
jccit@caltech.edu

Nicholas Adams-Cohen

Postdoctoral Scholar
Immigration Policy Lab
Stanford University
Stanford, CA 94305
nadamsco@stanford.edu

R. Michael Alvarez*

Professor of Political and Computational Social Science
Division of Humanities and Social Sciences
California Institute of Technology
Pasadena, CA 91105
rma@caltech.edu

November 17, 2020

Abstract

Social media data is now widely used by many academic researchers. However, long-term social media data collection projects, which most typically involve collecting data from public-use APIs, often encounter issues when relying on local-area network servers (LANs) to collect high-volume streaming social media data over long periods of time. In this technical report, we present a cloud-based data collection, pre-processing, and archiving infrastructure, and argue that this system mitigates or resolves the problems most typically encountered when running social media data collection projects on LANs at minimal cloud-computing costs. We show how this approach works in different cloud computing architectures, and how to adapt the method to collect streaming data from other social media platforms.

Keywords: Social media, Cloud computing, Twitter, Time series.

1 Introduction

Social media data is now widely used in many studies in computer and social science [1]. Many of these studies collect short-term cross-sectional samplings of social media data, while others take advantage of free-access social media APIs and attempt to build longer-term time series that monitor discussions and behaviors online. However, in any project that involves long-term social media data collection efforts, there are many potential issues with collecting a reliable and consistent pipeline of social media data when using on local-area network servers (LANs). In this technical report, we begin by discussing some of the issues that we have encountered in our own experience running a long-term Twitter data collection project (which at this point has been ongoing since 2014). We then present a cloud-based data collection, pre-processing, and archiving infrastructure which we argue mitigates or resolves many of the problems we have encountered, at minimal cloud-computing costs.

*We thank the John Randolph Haynes and Dora Haynes Foundation for supporting some of this research. We received support for our use of Google Cloud Platform through Google's COVID-19 research program. We also thank Anima Anandkumar and Anqi Liu for their work with us on related projects.

2 Problems collecting streaming social media

If a researcher is interested in quickly collecting cross-sectional social media data from Twitter, the use of the so-called “streaming” and “REST” APIs are relatively straightforward [2, 3]. Subject to rate limits, Twitter allows researchers to get access to a great deal of incoming Twitter data, including the content of a message, associated metadata, and information about the user account. In our application, where we are interested in studying the online conversations concerning different political and social topics, collecting data from the Twitter Streaming API by keyword or hashtag data filtering over a brief window of time is relatively straightforward, and is the methodology that many scholars use in their research [4, 5, 6, 7, 8].

This situation becomes more complicated if the research project involves longer-term monitoring of conversations and discussions on Twitter. For example, one of our ongoing projects involves monitoring Twitter mentions of voter issues during elections, requiring us to collect data continuously in the weeks before, during, and after an election. Ever since beginning this project in 2014, we have worked to refine and improve our methodology for collecting these data [9, 10]. Our process focuses on searching for specific keywords that are associated with topics including election fraud, voting by mail, and registering to vote. In another example of long-term social media monitoring, we are developing methods for collecting Twitter conversations using dynamic keyword selection in situations where the discussion might be rapidly-evolving over long periods of time [11].

In attempting to build long-term, multi-year, social media data collection projects on local machines, several prominent problems emerge. In our experience, simple issues can crop up. Our early work used Python scripts running on local university servers, connected to local-area networks. We found these scripts often encountered problems accessing the Twitter APIs, had trouble with network access, or competed with other processes running on the servers. Good programmers can often test scripts collecting longer-term social media data, identify some of these issues, and pause collection to revise code. [9].

However, even good programmers will have trouble resolving systems failures. An ideal social media monitor should maximize the amount of data gathered while minimizing the influence of any interruptions. Although a robust script that integrates diagnostic tools can help sustain the program, without solving the underlying system failures and limitations the threat of substantial data loss remains.

First, relying on local hardware introduces difficult, and in some cases, impossible-to-anticipate system failures. Power outages can knock systems offline, and without a secondary local system in place, data in a time series will be permanently lost. Network instability can also undermine data collection efforts, especially during peak-use hours or if network infrastructure is temporarily down for maintenance. Furthermore, if there is permanent system damage to a local system, it can be difficult, if not impossible, to recover data.

Second, the kinds of local systems most researchers have access to are not designed for the specific needs of collecting real-time streaming data. Collecting these data requires a system that can quickly and effectively obtain, buffer, process, and store large continuous streams of incoming information. Collecting this type of high-frequency, continuously streaming data involves specific computational considerations, with specialized algorithms designed to best capture these data [12]. Without a good system designed for these tasks, processing and saving files can temporarily interrupt the Twitter stream and limit the amount of data gathered. Slow processing that doesn’t meet the publishing speed will also be constantly disconnected from the social media streams. Given these interruptions are most likely to occur during periods of heavy Twitter traffic, the censored Twitter data may not be missing at random, with systemic omitted Twitter data from the streaming API potentially biasing the results of a study [13].

Finally, setting up a LAN can potentially limit future collection efforts. As a collection project naturally expands, computational power, active memory, and storage considerations may change. However, local systems can be difficult, expensive, and time-intensive to upgrade, especially if one needs to address these concerns repeatedly in a multi-year research effort. On the other hand, for seasonal projects such as election monitoring, local systems can be less efficient as they are difficult to downgrade or temporarily shutdown to cut expenses.

3 Cloud-based social media monitoring

In this paper, we present specific solutions to data collection on three popular cloud computing services: Google Cloud Platform (GCP), Amazon Web Services (AWS), and Oracle Cloud. While we illustrate our solution on these three platforms, the methods and processes we outline can be generally applied to other cloud services. While we do not claim here that we are the first to develop this type of workflow, we want to provide details about how our long-term social media data collection solution operates for other researchers to evaluate and utilize in their own work.

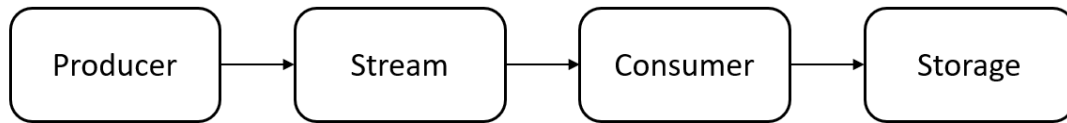


Figure 1: Social Media Monitor Workflow

Essentially we built a process that solves or mitigates many of the problems using LANs by moving the data collection and pre-processing steps onto cloud-computing platforms.²

In the next sections of this paper, we first describe the general workflow of our social media monitoring system before providing detailed guidelines on how the process works in GCP, AWS, and Oracle Cloud.

3.1 Workflow

We developed an architecture using cloud resources to tackle the problems discussed in the previous section. This system collects as much social media data as rate limits allow in a stable and failure-tolerant manner.

The system consists of four parts: a data producer, a data stream, a data consumer, and storage. As shown in Figure 1, it starts with a producer that requests social media data from the API and acts as a data provider to the other parts of the system. Then the produced data are published in the data stream for temporary storage. The data stream secures every published record on a timeline according to the timestamp. Before the records expire, a data consumer retrieves them from the data stream and either sends them to analytics modules or to the short-term/long-term storage solution.

We describe the workflow of Twitter monitoring step-by-step in the following subsections and will expand the discussion to other social media platforms in the next section.

3.2 Data Producer

The data producer runs Python scripts on the cloud compute instances and connects two parts of our data pipeline: the Twitter Stream API and the cloud data stream.

The data producer:

- Uses the `TwitterAPI` package to access the Twitter Stream API, `google-cloud` to interact with GCP, `boto3` to interact with AWS, and `oci` to interact with Oracle Cloud.
- Connects to the Twitter Stream API using Twitter Developer credentials.
- Connects to the cloud data stream using credentials/tokens.
- Requests tweets from the Twitter Stream API.
- Publishes streaming tweets one-by-one to the cloud data stream.

We use the Twitter Stream API to collect real-time tweets. The most commonly used Twitter APIs for collecting tweets are the Stream and REST APIs. The Twitter Stream API delivers real-time tweets continuously as soon as they are published and prior to subsequent alterations, such as deletion or censoring by the platform. Furthermore, as tweets are collected right after they are published, they do not carry information regarding the retweets and quotes that occur after a message is sent. On the contrary, the Twitter REST API allows users to search for the tweets that were posted in the past 7 days (30-day and full-archive endpoints are available upon upgrade). These tweets contain the latest retweet and quote information up to the time of extraction. However, instead of delivering all tweets that meet the filtering rules, the REST API only returns a subset of them that are most relevant. Therefore, to collect maximal tweets, our system focuses on the Stream API, only relying on the REST API as a backup in the case of interruptions.

It is important to note that the rate limits of the Twitter Stream API can cause potential data loss. The Twitter Stream API has a rate limit that only allows the delivery of up to 50 tweets/second. Some projects such as COVID-19 have exclusive unlimited endpoints. Those unlimited endpoints should be prioritized as they have no rate limits.

²To reduce cloud-computing data storage costs, and to make the social media data we collect more readily accessible to our research group, we outline a process of piping the pre-processed data to cheap, secure, and easy-to-use data storage applications (here Google Drive). Note that, budget allowing, all data can be stored on a single platform.

Selecting a suitable cloud compute instance is also important, as it ensures sufficient computing power at the lowest price. In our COVID-19 project, we use two data producers to request tweets of independent topics. We found that a Python script that produces at Twitter’s rate limit (50 tweets/second) occupies on average 20% of a 2GHz CPU core, and running data producers on separate CPU cores can prevent the processes from competing with each other for resources. For a two-producer project like ours, we suggest using an instance with two CPU cores. We further recommend allocating more than 1GB of active memory space to compile large packages and run data consumer scripts. Table 1 shows three popular compute instances in GCP, AWS, and Oracle Cloud that each has specifications that meet the base requirements of the two-producer project. While the AWS system costs less than the others, AWS employs a baseline performance mechanism³ that limits the instance’s performance once the average CPU usage exceeds the baseline. For the t3.medium instance, performance will be significantly reduced if the user constantly uses more than 20% of the CPU(s). In general, the performance/price ratios are similar across platforms. Users can always customize or resize their instances to balance the performance and cost.

Table 1: Comparison of Instance Specifications

| Specs | GCP Compute VM n2-custom | AWS EC2 EC2 t3.medium | Oracle Compute VM.Standard.E2.2 |
|--|--------------------------|-----------------------|---------------------------------|
| CPU Series | Intel(R) Xeon(R) | Intel(R) Xeon(R) | AMD EPYC |
| CPU Model | | Platinum 8175M | 7551 |
| Core Clock | 2.80GHz | 2.50GHz | 2.00GHz |
| Core Count | 2 Cores | 2 Cores | 2 Cores |
| Memory | 4 GB | 4 GB | 15 GB |
| Baseline Performance (per CPU)* | 100% | 20% | 100% |
| Pricing | \$50/Month | \$30/Month | \$61/Month |

* CPU performance is restricted if average CPU usage exceeds the baseline.

3.3 Data Stream

Cloud data streams are temporary storage services designed specifically for streaming data. Their role in the data collection pipeline is like librarians – they receive, organize, and preserve new information from the source and provide multiple means for the users to access the collected information. Data streams are the most essential part of the social media monitor, as they handle heavy data traffic that is beyond most local systems’ capabilities, and make each streaming record available to all parts of the system. Data streams improve the reliability of the whole process by making it robust to failures caused by peaks of incoming data and serving as buffers for the subsequent parts of the workflow to re-visit data records if an error should occur.

While the data stream services have different names and specifications in GCP, AWS, and Oracle Cloud, they serve for the same purpose and work in the same way. In the following discussion, we talk about the rate limits, retention periods, and pricing methods of these stream services.

Table 2 shows the default rate limits for the stream services Pub/Sub (GCP), Kinesis (AWS), and Oracle Stream. In the three services, Kinesis and Oracle Stream operate on basic units (called shards in Kinesis, and partitions in Oracle Stream), while Pub/Sub runs as a whole stream.

The basic units are pre-set components of a stream service. They are restricted by default rate limits and are independent of each other. Inside these units, published records are ordered by the timestamps of the “put” events. A shard of the Kinesis stream supports up to 1,000 records/s or 1 MB/s (whichever is met first) for writes, and 5 requests/s or 2 MB/s for reads. For Oracle Cloud, a partition has similar rate limits to a Kinesis shard, except it has no restrictions on the number of write requests per second.

Users can create more units in their stream services to meet increased demand. In AWS Kinesis, a Python script can use metrics from AWS Cloudwatch to monitor the usage of Kinesis shards, and scale up the number of shards if thresholds are met. For example, to prevent the incoming tweets from hitting the writing rate limits and causing data loss, we can use upper bounds 800 KB/s and 800 records/s as signals for immediate shard-creation. Once the average incoming tweets exceed one of the thresholds, the Python script immediately creates a new shard and lowers the average burden. We can also set lower bounds 500 KB/s and 500 records/s as signals for delayed shard-deletion. Once the average records per shard fall below the lower bounds, we want to delete one or more shards to reduce cost. To be conservative

³<https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/burstable-credits-baseline-concepts.html>

Table 2: Comparison of Cloud Data Streams*

| Specs | | GCP Pub/Sub | AWS Kinesis | Oracle Stream |
|------------------|-------|---|---------------------------|----------------------------|
| Unit Limits | Write | None | 1 MB/s 1,000 Records/s | 1 MB/s Unlimited Writes |
| | Read | | 2 MB/s 5 Reads/s | 2 MB/s 5 Reads/s |
| Project Limits | Write | 50 MB/s (small**) 200 MB/s (large) | Unit Limit*N | Unit Limit*N |
| | Read | 100 MB/s (small) 400 MB/s (large) | Unit Limit*N | Unit Limit*N |
| Retention Period | | 7 Days | 24 Hours | 24 Hours |
| Pricing | | Message Ingestion, Delivery, and Storage | Shards*Hour | Partitions*Hour |

* Default rate limits can be increased upon upgrade.

** Large regions: europe-west1, us-central1, us-east1. Small regions: other regions.

in deleting shards, we recommend waiting at least three hours after the lower bounds signals were triggered. The records that were published in the deleted shards will be available until expiration. For Oracle Cloud, it is not possible to add new partitions to existing streams, so users need to plan ahead or manually migrate to a new stream that has more partitions if incoming record surges.

The GCP Pub/Sub does not involve stream units explicitly. Instead of focusing on unit limits, Pub/Sub users must pay attention to project and resource limits. For a project that is located in large regions (europe-west1, us-central1, us-east1), all Pub/Sub topics combined cannot exceed 200 MB/s for writes and 400 MB/s for reads. The corresponding rate limits in other regions are 50 MB/s and 100 MB/s. The users should create separate GCP projects for large social media monitors that can potentially exceed the project limits.

The retention period determines how long a record is available to be read after being published in the stream. The default retention periods for Kinesis and Oracle Stream are both 24 hours, and it is 7 days for Pub/Sub. Users can upgrade to a longer retention period when creating the stream service.

There are two pricing methods, fixed pricing per unit*hour and flexible pricing per volume of data transmission and storage. Kinesis and Oracle Stream use the former and Pub/Sub uses the latter. The fixed pricing method leaves it to users to optimize the size and cost of the stream, and there are inevitable resource losses in running the units below the rate limits and in re-sizing the stream. On the contrary, flexible pricing only charges the resources being used. However, we wish to emphasize that flexible pricing does not definitively lead to a lower overall cost.⁴

3.4 Data Consumer

Data consumers are cloud services or customized programs that read records from the data streams and send them to either cloud/local databases for storage or to analytical modules for data processing and analyses. In our social media monitors, we use DataFlow in GCP, Firehose in AWS, and the "get_messages" function in Oracle Cloud to extract records from the data stream. We prefer using a cloud service like DataFlow and Firehose instead of relying on customized programs given these services tend to have lower latency and higher stability. To add one more layer of safety, we send the extracted data immediately to cloud databases for short term storage and subsequently invoke BigQuery (GCP), Lambda Functions (AWS), or Data Analytics (Oracle Cloud) for analyses.

3.5 Storage

Our social media monitoring system puts collected data temporarily in cloud storage before archiving data in Google Drive folders. Cloud storage services, such as Cloud Storage (GCP), S3 (AWS), Object Storage (Oracle Cloud) are natural data transfer and storage solutions between stream services and compute instances. They are ideal for frequently used data but are not cost-efficient for large sets of raw social media data collected over a long period of time. After the data collection pipeline is checked and the real-time analyses are completed, we store the data summaries and results of the analyses in a cloud MariaDB database and transfer the raw data to a Google Drive shared with all members of the research team. This three-layer storage structure (cloud data stream–cloud storage–Google Drive) is robust to system

⁴Please refer to the pricing pages and the cost calculators for each cloud service to better estimate prices.

failures given it inherits the stability and compatibility from the cloud services, and, if an error should happen, the three-layered system can easily recover data from a previous step.

4 Other Social Media Platforms

Of course, Twitter is not the only source of dynamic social media data for researchers. While the Twitter API's relative ease of use and open policy make it one of the most popular sources of data in academic studies, many researchers conduct research utilizing data from other popular social media platforms with APIs, such as Reddit [14, 15], YouTube [16, 17], and Facebook [18, 19].

Our approach can be adapted to collect streaming or dynamic social media data from these other platforms. Considering Figure 1, the main modification in our process would be in the **Producer** step of the workflow. In this step, we describe setting up a virtual compute instance capable of running code designed to interact with the Twitter API. To modify our process to collect other forms of streaming data, we would simply rewrite this code to interact with another platform's API.

For example, suppose we wished to collect data over time from a particular Reddit community (a subreddit). We could reuse the same process described in the previous section, revising the set of scripts in the Producer step to access data from the Reddit API. There might also be some minor changes required in the code used in the preprocessing and consumer steps, as the data stream from the Reddit API may differ from the Twitter API. These minor alterations aside, conceptually our approach is highly adaptable for collecting streaming and dynamic data from a variety of social media platforms, as long as they have a public API or their data can easily be obtained via a script that can run on a cloud computing instance.

5 Discussion and Conclusion

Many research groups are using social media data in their studies of political, social, and economic attitudes and behavior. Interest is increasing in the development of longer-term datasets that can be used to analyze changes over time in attitudes and behavior [8]. However, in our efforts to collect longer-term social media datasets from local servers using local-area networks, we have encountered important limitations in the availability and reliability of those systems for these purposes.

To improve the reliability of our longer-term social media data collections process, we have developed a cloud-based infrastructure that is adaptable to several platforms. By moving the data collection and pre-processing stages into the cloud, we avoid many of the problems encountered when relying on local servers and LANs. We designed our process to make data easily accessible to our research group by moving the data to a secure and usable storage solution like Google Drive. We have shown that these processes work across different cloud computing systems, and can be used to collect both Twitter and other streaming social media data.

There are many opportunities for researchers who are interested in using social media data to study the longer-term dynamics of political, social, and economic attitudes and behavior. We hope that by providing the details of our framework, other researchers can evaluate and potentially use this report as initial guidance in adopting a cloud-based process, improving their ability to collect similar datasets.

References

- [1] Marko Klasnja and Pablo Barbara and Nicholas Beauchamp and Jonathan Nagler and Josh Tucker. Measuring public opinion with social media data. In *The Oxford Handbook of Polling and Survey Methods*, edited by Lonna Rae Atkeson and R. Michael Alvarez, pages 555–582. Oxford University Press, 2018.
- [2] Matthew A. Russell. *Mining the Social Web: Analyzing Data from Facebook, Twitter, LinkedIn, and Other Social Media Sites*. O'Reilly Media, Inc., Second Edition, 2014
- [3] Zachary Steinert-Threlkeld *Twitter as Data*. Elements in Quantitative and Computational Methods for the Social Sciences. Cambridge University Press, 2018
- [4] Brendan O'Connor, Ramnath Balasubramanyan, and Bryan R. Routledge. From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media, 2010
- [5] Michael D. Conover, Bruno Gonçalves, Alessandro Flammini and Filippo Menczer. Partisan asymmetries in online political activity. *EPJ Data Science* 1 (6), 2012.

- [6] Pablo Barberá and Gonzalo Rivero Understanding the political representativeness of Twitter users. *Social Science Computer Review*, 33 (6), 2014.
- [7] Nicholas Beauchamp. Predicting and Interpolating State-Level Polls Using Twitter Textual Data. *American Journal of Political Science*. 61 (2), 2017
- [8] Nicholas Adams-Cohen. Policy Change and Public Opinion: Measuring Shifting Political Sentiment with Social Media Data. *American Politics Research*, forthcoming 2020.
- [9] Nicholas J. Adams-Cohen and Clare Hao and Cherie Jia and Nailen Matschke and R. Michael Alvarez. Election Monitoring Using Twitter. Caltech/MIT Voting Technology Project Working Paper 129, 2017. vote.caltech.edu/working-papers/129
- [10] R. Michael Alvarez, Nicholas Adams-Cohen, Seo-young Silvia Kim, and Yimeng Li. *Securing American Elections: How Data-Driven Election Monitoring Can Improve Our Democracy*. Elements in Campaigns and Elections. Cambridge University Press, <https://doi.org/10.1017/9781108887359>.
- [11] Anqi Liu and Maya Srikanth and Nicholas Adams-Cohen and R. Michael Alvarez and Anima Anandkumar. Finding social media trolls: Dynamic keyword selection methods for rapidly-evolving online debates. AI For Social Good Workshop, NeurIPS 2019. [arXiv:1911.05332\[cs.LG\]](https://arxiv.org/abs/1911.05332)
- [12] Brian Babcock, Shivnath Babu, Mayur Datar, Rajeev Motwani, and Jennifer Widom. Models and issues in data stream systems. SIGMOD/PODS02: International Conference on Management of Data and Symposium on Principles Database and Systems, 2002.
- [13] Fred Morstatter, Jürgen Pfeffer and Huan Liu. When is it biased?: assessing the representativeness of Twitter’s streaming API. Web Science Track at WWW, 2014.
- [14] Chenhao Tan, Vlad Niculae, Cristian Danescu-Niculescu-Mizil, and Lillian Lee. Winning arguments: Interaction dynamics and persuasion strategies in good-faith online discussions. Proceedings of the 25th international conference on world wide web, 613–624, 2016.
- [15] Rishab Nithyanand, Brian Schaffner, and Phillipa Gill. Online Political Discourse in the Trump Era [arXiv:1911.05332\[cs.CY\]](https://arxiv.org/abs/1911.05332), 2016.
- [16] Jean Burgess and Ariadna Matamoros-Fernández. Mapping sociocultural controversies across digital media platforms: one week of #gamergate on Twitter, YouTube, and Tumblr Communications Research and Practice, 2(1). 2016.
- [17] Bernhard Rieder, Ariadna Matamoros-Fernández, Óscar Coromina From ranking algorithms to ‘ranking cultures’: Investigating the modulation of visibility in YouTube search results *Convergence*, 24(1) 50–68. 2016.
- [18] Bernhard Rieder, Rasha Abdulla, Thomas Poell, Robbert Woltering, and Liesbeth Zack. Data critique and analytical opportunities for very large Facebook Pages: Lessons learned from exploring “We are all Khaled Said” *Big Data & Society*, <https://doi.org/10.1177/2053951715614980>. 2015.
- [19] Bente Kalsnes. The Social Media Paradox Explained: Comparing Political Parties’ Facebook Strategy Versus Practice. *Social Media + Society*, April 2016. [doi:10.1177/2056305116644616](https://doi.org/10.1177/2056305116644616).