# Detecting Twitter Bots That Share SoundCloud Tracks

**Axel Bruns, Brenda Moon, Felix Victor Münch, Patrik Wikström**
Digital Media Research Centre
Queensland University of Technology
a.bruns@qut.edu.au, brenda.moon@qut.edu.au, felixvictor.muench@hdr.qut.edu.au, patrik.wikstrom@qut.edu.au

**Stefan Stieglitz, Florian Brachten, Björn Ross**
Research Group Professional Communication in Electronic Media / Social Media
University of Duisburg-Essen
stefan.stieglitz@uni-due.de, florian.brachten@uni-due.de, bjoern.ross@uni-due.de

## ABSTRACT

Sharing platforms for creative content are often closely connected to general purpose social media platforms like Twitter. This also means that coordinated and automated mechanisms for promoting such content are likely to span both sites: spammers and bots operate across both platforms. This work-in-progress paper presents first results from an effort to develop activity metrics that enable the detection of Twitter bots promoting SoundCloud content.

## CCS CONCEPTS

• **Information systems → World Wide Web → Web applications** → Social networks

## KEYWORDS

SoundCloud, Twitter, automated activity, bot detection

## 1 INTRODUCTION

The presence of bots on Twitter has come under increasing scrutiny in recent years; some studies estimate that bots represent up to 15% of Twitter accounts [1]. Such bots may be entirely automated, or guided at least in part by human controllers; they might perform legitimate automatic tasks (e.g. posting regular weather updates or alerting their followers to new earthquake activity), or engage in more subversive activity designed to promote certain political actors and ideologies or to help specific keywords and hashtags achieve "trending topic" status. In the latter cases, their nature as automated bots is often concealed more or less carefully, in order to trick ordinary human users of Twitter into believing that the bot-promoted content they encounter in their timelines or trending topics lists is prominent as the result of organic user activity on the platform, rather than because of a concerted effort to game the platform's built-in affordances and algorithms. Literature that categorises social bots and discusses their positive and negative impact as well as their ethical implications remains limited due to the comparative recency of the phenomenon [2, 3, 4].

Much of the research into the latter category of bots has focussed on their uses for political purposes (e.g. [5]); this has been especially prominent in the context of concerns about the deliberate promotion of "fake news" and other forms of propaganda, misinformation, and disinformation by various political and state actors [6]. By contrast, bot activity in commercial contexts has received considerably less attention to date. To address this lack of knowledge, the present paper focusses on the commercial and quasi-commercial uses of such bots: it investigates the use of Twitter bots to promote content from the popular audio-sharing site SoundCloud, and proposes a number of social media metrics that may be used to detect bot-like behaviour in the sharing of such content. Such metrics may then also be utilised in other bot detection contexts where bots are suspected to provide an artificial boost to the visibility of content by sharing its URLs on Twitter – including also the sharing of "fake news" URLs and related political content.

SoundCloud provides a useful test case for this approach. In spite of recent concerns about its long-term financial viability, it remains a popular site for the sharing of music by professional and amateur musicians, as well as of other audio content including journalistic interviews, podcasts, radio shows, and other formats. Such content is then also widely shared through social media; our research finds

that, on average, SoundCloud URLs were shared on Twitter more than 175,000 times per day over the past year. This provides a rich dataset for the detection of bots that seek to boost the visibility of SoundCloud content on Twitter.

Further, simple Web searches show the ready availability of commercial services that offer various packages for SoundCloud and Twitter promotion: of these, some vendors sell native SoundCloud plays, likes, comments, followers, and downloads that are designed to boost the on-site metrics of a given track and thereby enhance its visibility on SoundCloud itself (Fig. 1), while others offer networks of Twitter bots that – amongst other purposes – can be used to mass-promote SoundCloud content by posting hundreds and thousands of original tweets in parallel, or by retweeting a given original post that links to a SoundCloud URL. Given the existence of this ecosystem for the automated promotion of audio content on SoundCloud itself as well as on Twitter, an analysis of promotional practices on either platform should find ample evidence of bot activity.



**Figure 1: A typical SoundCloud metrics vendor's page.**

As part of a larger research project that investigates patterns of user activity on SoundCloud, this paper focusses on the role of Twitter bots in promoting SoundCloud content and develops various metrics that facilitate their detection; a related paper, using SoundCloud data, examines native promotional activity on the site itself [7]. Further research will correlate the findings from SoundCloud and Twitter in order to detect cross-platform patterns of bot activity.

## 2 APPROACH

### 2.1 Datasets
We build on a longitudinal dataset of tweets (and their associated metadata) that shared SoundCloud URLs, gathered by regularly querying the public Twitter Search API for tweets that contained "soundcloud.com" (both verbatim within the tweet texts themselves, and as the eventual destination of URLs shortened by Twitter's built-in URL shortener *t.co*). The tracking of such tweets since 2 February 2017 had resulted in a dataset of more than 60 million such tweets by early January 2018; from these, to arrive at a manageable dataset we have selected all tweets that were posted during the months of March and April 2017 for our further analysis. This leaves 11,530,680 tweets, posted by 2,099,526 unique Twitter accounts.

We further processed the SoundCloud URLs contained in this dataset in order to strip out any extraneous modifiers (such as "?utm_medium=twitter" and other URL modifiers employed by Google Analytics) and systematise the formatting of these URLs (for instance reformatting all *m.soundcloud.com* URLs, as posted from mobile devices, to standard *soundcloud.com* URLs). This ensures our ability to aggregate all of the tweets that share a given SoundCloud content item to a single count.

From this processed dataset, we extracted two key subsets for further analysis. First, for all reformatted URLs that followed the standard pattern of *soundcloud.com/[user name]/[track name]*, we generated a count of the total number of tweets that shared these tracks during the two months, and selected those tracks that were shared on Twitter at least 1,000 times during these two months (leaving 223 distinct SoundCloud tracks that were shared in 914,131 tweets, out of 2.1m distinct tracks shared during the two months). This constitutes our Top Tracks dataset.

Second, we also generated a count of the number of times that each Twitter account in our dataset had shared any SoundCloud track, and selected those accounts that had shared SoundCloud tracks in at least 500 tweets (including retweets) during the two months (leaving 649 distinct accounts that posted over 1m tweets sharing SoundCloud tracks). This constitutes our Top Accounts dataset.

We concentrate on these top tracks and top accounts because – assuming that bot-based promotion is effective in the first place – it is in this most visible group of tracks that we most expect to be able to detect bot-like promotional patterns, and because we assume that bot accounts that are used to promote SoundCloud content will share such content more often than most ordinary users.

### 2.2 Key Metrics
For these Top Tracks and Top Accounts datasets, then, we devised several metrics that appear useful in the detection of bot-like behaviours. We define the following measures:

1. $Account\ Diversity = \frac{\#\ Unique\ Twitter\ Accounts}{\#\ Tweets}$

Calculated for each SoundCloud track, this value may range from close to 0 to 1: if a given SoundCloud track was shared by only one Twitter account, Account Diversity would calculate as $\frac{1}{n}$ (for $n \geq 1000$, given the filter we have applied already), while if each of $n$ tweets sharing the track was posted by a different account, then the Account Diversity value would be $\frac{n}{n} = 1$. Low Account Diversity values thus mean persistent promotion by a small number of accounts (and indicate possible spamming, by bots or humans), while high Account Diversity values might appear more organic but could still point to concerted promotional efforts by a network of bots.

2. $Tweet\ Originality = \frac{\#\ Non\text{-}Retweets - \#\ Retweets}{\#\ Tweets}$

This metric measures the relative contribution of retweets and non-retweets to a given case. To aid comparison, we normalise its value to between -1 and +1: calculated for each SoundCloud track, a Tweet Originality value of -1 would mean that all tweets sharing the track were retweets, while a value of +1 would mean that each of the tweets was an original tweet or @mention. The former would point to considerable resharing of one or more initial tweets linking to the track; the latter may indicate a concerted campaign to post many original tweets, and/or to @mention a large number of accounts while sharing the track, and both such practices might be considered a form of spam if they are undertaken at high volume by a small number of accounts.

Further, Tweet Originality may also be calculated across the tweets posted by each Twitter account, independent of which SoundCloud tracks they link to. Here, a value of -1 would mean that the account has only ever retweeted posts that link to such tracks; in doing so it amplifies the messages of other accounts but does not contribute any original posts of its own. Conversely, a value of +1 means that each of the account's tweets that link to a SoundCloud track has been individually crafted. Neither of these behaviours is in itself inherently more indicative of bot-like activity (human users may be genuinely committed to sharing newly found SoundCloud tracks with their followers, or to on-sharing what their Twitter sources have found), but may appear more suspicious in combination with unusual patterns across the other metrics introduced here.

3. $Tweet\ Text\ Diversity = \frac{\#\ Unique\ Tweet\ Texts}{\#\ Tweets}$

Again, if calculated for each SoundCloud track this value ranges from close to 0 to 1: if all of the tweets promoting a SoundCloud track were identical, the Tweet Text Diversity value would be $\frac{1}{n}$; if they were each different, it would be $\frac{n}{n} = 1$. Low Tweet Text Diversity values would mean that the track is being promoted on Twitter in a highly uniform way; this could point to the prominence of one widely retweeted post, if most of these tweets are retweets. Alternatively, in the absence of substantial retweeting a low Tweet Text Diversity would mean that there are many apparently independent tweets with identical content; this could be the result of an organic promotion campaign (e.g. through an artist Website that enables visitors to post a prepared tweet to express their fandom) or of coordinated bot activity.

Additionally, Tweet Text Diversity may also be calculated across the tweets posted by a given Twitter account. Here, low Tweet Text Diversity would mean the repeated posting of identical tweets (which would usually indicate a form of spamming), while high Tweet Text Diversity should be more common for most Twitter accounts, which are unlikely to post the exact same tweets repeatedly.

(In the present study, we take a simple approach to determining the number of unique tweets in the dataset, by comparing the entire text of each tweet against others. A further extension of this approach could be to allow for fuzzy matching, to treat tweets that are *almost* entirely identical to each other as instances of the same message. This could account for spamming practices where random characters are added to tweets in order to fool Twitter's spam detection algorithms; for the present work-in-progress paper, however, we have not implemented such matching.)

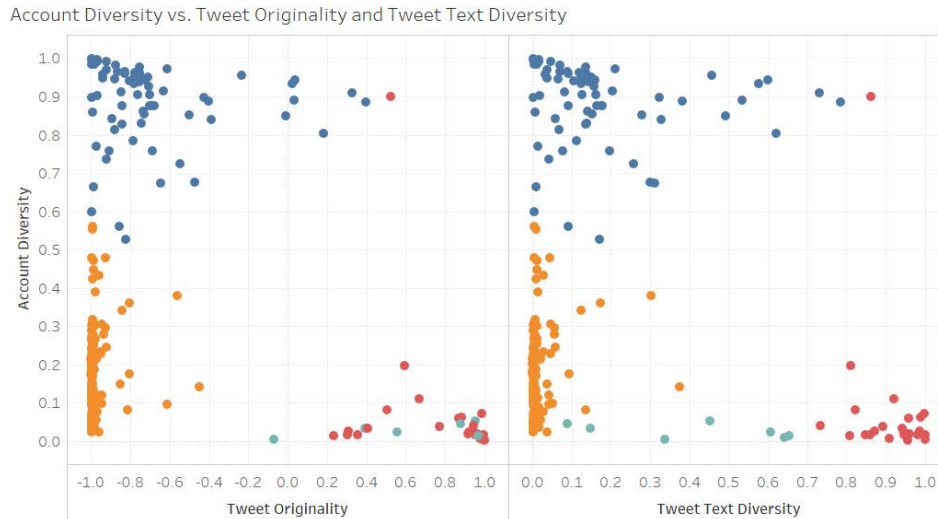4. $SoundCloud\ URL\ Diversity = \frac{\#\ Tracks\ Shared}{\#\ Tweets}$

This measure is calculated for each Twitter account that shared SoundCloud tracks, and again results in values from close to 0 to 1: if each tweet posted by the account shared the same track, it would equate to $\frac{1}{n}$ (with $n \geq 500$, given the cutoff we used to identify Twitter accounts in our dataset); if each tweet shared a different track, then the SoundCloud URL Diversity measure for the account would be $\frac{n}{n} = 1$. Neither of these points in itself to more bot-like behaviour: while excessive focus on a single track should certainly be regarded as suspicious, a high diversity of interests coupled with a substantial volume of tweets could point to a bot that is programmed to post links to a long list of tracks. This measure – and indeed all four of the measures we have introduced here – are therefore useful especially if they are deployed in combination with each other, and with other more basic metrics.

## 3 PRELIMINARY ANALYSIS

In the remainder of this paper we present the early results from our implementation of these metrics on the two datasets. As work in progress, our findings should not yet be regarded as conclusive in their own right, but point to a number of major activity patterns that offer potential for further, more detailed analysis; such analysis may involve a larger subset of our entire, year-long dataset and employ less restrictive cutoffs for the top tracks and accounts, and should also include a more extensive qualitative review of the key SoundCloud tracks and Twitter accounts revealed by the quantitative analysis. These steps are left for further work on this project. Here, we examine instead the interactions between the different metrics already available, and highlight key areas for further investigation.

### 3.1 Metrics per SoundCloud Track

We begin by examining the distribution of these metrics across the SoundCloud tracks in our Top Tracks dataset that were shared 1,000 times or more during March and April 2017 (Fig. 2). Here and in section 3.2, for the purposes of a preliminary distinction between tracks that exhibit a similar clustering of values across the three metrics we draw on k-means clustering, as implemented in the data

**Figure 2: Account Diversity compared to Tweet Originality and Tweet Text Diversity (for tracks shared ≥1,000 times); tracks grouped by clustering across the three metrics, using k-means clustering in *Tableau*.**

visualisation package *Tableau*. This results in a number of distinct subsets, shown in Figs. 2 and 3 in different colours.

It is notable, first, that Tweet Originality and Tweet Text Diversity appear broadly correlated, and this is to be expected: if a track is mainly shared in a large number of retweets, it appears likely that these result from the widespread on-sharing of a small number of prominent original tweets, rather than from a multitude of low-volume retweeting of a larger number of original tweets. Indeed, Fig. 2 shows that the vast majority of our 223 tracks with 1,000 or more shares during March and April 2017 (184 tracks, or 83%) exhibit Tweet Originality values below 0: they were mainly shared through retweeting. For most of these tracks, retweets far surpass non-retweets, in fact: 150 tracks reach Tweet Originality scores of 0.8 or above. These tracks further subdivide along Account Diversity lines, however: one group of 72 tracks (shown in blue) was retweeted by a highly diverse range of accounts (Account Diversity ≥ 0.5), while another group of 120 tracks (in orange) appears to have been retweeted repeatedly by a smaller number of accounts; this raises suspicions of bot-like activity which must be investigated further.

By contrast, the 31 tracks shown in red and turquoise in fig. 2 are distinguished by their high Tweet Originality, indicating a relative absence of retweets, combined with a low Account Diversity – this means that a small number of Twitter accounts repeatedly shared each track in multiple original tweets. Further, this group subdivides along Tweet Text Diversity lines: 24 of these tracks (in red) show a high Tweet Text Diversity of 0.7 and above (as their promoters shared them time and again, they varied the language they used in doing so), while 7 show much lower Tweet Text Diversity (the same tweets were posted over and over again). Taken to their extreme, both practices may be considered as spamming; of these, the former practice may be more actively seeking to evade any penalties for doing so,

by regularly varying the message associated with the SoundCloud link and thereby appearing less repetitive. It seems likely that some of the accounts involved here will represent the tracks' artists or record labels, and we intend to examine these connections in further analysis.
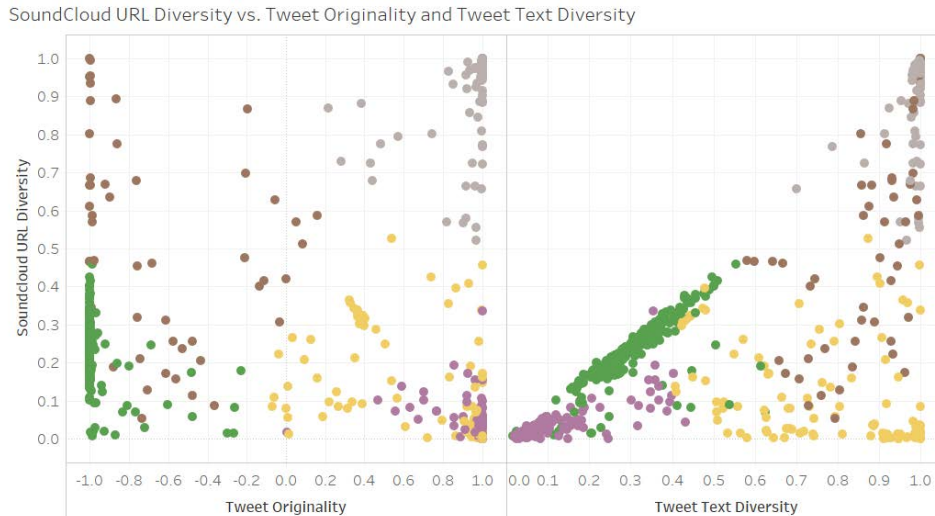
### 3.2 Metrics per Twitter Account

The distribution of metrics in our Top Accounts dataset, containing the Twitter accounts that shared links to SoundCloud tracks at least 500 times during March and April 2017 (fig. 3), is less distinct, but nonetheless points to a number of diverging patterns that require further analysis. Again, we use k-means clustering to distinguish subsets of accounts with divergent combinations of metrics.

First, it is evident that the Tweet Originality values across this account population are highly polarised. Of the 649 distinct accounts, 305 (47%) have a Tweet Originality score of -0.6 or less, and thus mainly engage in retweeting other accounts' links to SoundCloud tracks; another 270 (42%) achieve a score of 0.6 or above, and are mainly posting original tweets and @mentions that link to such tracks. Accounts with a balanced tweeting behaviour are far more rare: only 74 (11%) of the most active sharers combine original recommendations and retweets in a more even mix.

The group of heavy retweeters further subdivides into those accounts that share a wide range of tracks through such retweets (largely shown in brown), and those that focus mainly on promoting the same tracks over and over (mostly shown in green). Drawing on the k-means cluster analysis for our preliminary examination of these patterns, our distinctions between these groups are further complicated by the Tweet Text Diversity scores, however: while for the group of the heavy retweeters of a limited number of tracks (278 accounts, shown in green) the Tweet Text Diversity grows largely in direct relation with the SoundCloud URL Diversity, the other, smaller group (51

**Figure 3: SoundCloud URL Diversity compared to Tweet Originality and Tweet Text Diversity (accounts with ≥500 tweets); accounts grouped by clustering across the three metrics, using k-means clustering in *Tableau*.**

accounts, in brown) shows high Tweet Text Diversity independent of the range of SoundCloud tracks they promote.

Conversely, amongst the accounts with high Tweet Originality Scores we distinguish three broad groups. The first (83 accounts, in grey) combines its preference for original tweets with high scores for Tweet Text Diversity and SoundCloud URL Diversity: in other words, almost every tweet shares a different track, using different text. Such variation may indicate of strong personal commitment to SoundCloud, or result from automatic mechanisms that post a link to each new track appearing on SoundCloud.

Two other groups exhibit much lower SoundCloud URL Diversity. One of these (112 accounts, in yellow) maintains comparatively high Tweet Text Diversity: in other words, these accounts share a small number of tracks in ever-changing original tweets. By contrast, the other (125 accounts, in purple) shows very low Tweet Text Diversity: here, the same tracks are shared over and over again, largely using the same texts, but without resorting to retweets (and this practice may again be considered as a form of spam). Here, too, some such behaviours may point to the work of artists or labels in promoting their latest releases; in further work, we intend to investigate whether such monotonous promotion targets those tracks that we identified as benefitting from potentially spammy promotional activity in 3.1 above.

## 4   NEXT STEPS

This preliminary analysis of a two-month subset of our larger dataset points to a number of divergent patterns in how SoundCloud tracks are promoted on Twitter. Further qualitative exploration of the various groups of tracks and accounts we have identified, and additional correlation between per-track and per-account metrics, will yield further insight into these patterns, and enable the fine-tuning and further development of sharing metrics. This should also generate detection thresholds for these metrics, individually

and in combination, beyond which the likelihood of botness increases. Building on the observations from our initial two-month pilot study, we then intend to apply these heuristics to the larger, live dataset we are collecting, to examine sharing behaviour in action. Through this further work, we hope to identify bots and similar automated accounts more effectively. We also intend to examine activity patterns on the SoundCloud platform itself [7], to test whether the tracks found to be artificially hyped on Twitter also exhibit unusual activity patterns on SoundCloud.

## REFERENCES

[1] O. Varol, E. Ferrara, C.A. Davis, F. Menczer, and A. Flammini. 2017. Online Human-Bot Interactions: Detection, Estimation, and Characterization. *arXiv:1703.03107v2 [cs.SI]*. https://arxiv.org/abs/1703.03107

[2] T. Bucher. 2014. About a Bot: Hoax, Fake, Performance Art. *M/C Journal* 17, 3 (2014). http://journal.media-culture.org.au/index.php/mcjournal/article/view/814

[3] B. Moon. 2017. Identifying Bots in the Australian Twittersphere. In *Proceedings of the 8th Conference on Social Media & Society*.

[4] S. Stieglitz, F. Brachten, B. Ross, and A. Jung. 2017. Do Social Bots Dream of Electric Sheep?. In *Proceedings of the 28th Australasian Conference on Information Systems (ACIS)*.

[5] F. Brachten, S. Stieglitz, L. Hofeditz, K. Kloppenburg, and A. Reimann. 2017. Strategies and Influence of Social Bots in a 2017 German State Election. In *Proceedings of the 28th Australasian Conference on Information Systems (ACIS)*.

[6] C. Shao, G.L. Ciampaglia, O. Varol, A. Flammini, and F. Menczer. 2017. The Spread of Fake News by Social Bots. *arXiv:1707.07592*. https://arxiv.org/abs/1707.07592

[7] B. Ross, F. Brachten, S. Stieglitz, P. Wikström, B. Moon, F. Münch, and A. Bruns. 2018. Social Bots in a Commercial Context – A Case Study on SoundCloud. Paper presented at *European Conference on Information Systems (ECIS)*, Portsmouth, 23-28 July 2018.