

# A Tweet Consumer's Look At Twitter Through Linked Data Goggles Via Google Analytics

Thomas Steiner and Arnaud Brousseau

Google Germany GmbH, ABC-Str. 19, 20354 Hamburg, Germany  
`{tomac, arnaudb}@google.com`

**Abstract.** Twitter Trends<sup>1</sup> allows for a global or local view on “what’s happening in my world right now” from a tweet producers’ point of view. In this paper, we discuss the possibility to complete the functionality provided by Twitter Trends by having a closer look at the other side: the tweet consumers’ – i.e., readers’ – point of view. While Twitter Trends works by analyzing the frequency of terms and their velocity of appearance in tweets being written, our approach is based on the popularity of extracted named entities (in the sense of Linked Data) in tweets being read. Our experimentation architecture takes advantage of the possibility to use a client-side browser extension to harvest and dissect tweets from users’ timelines, i.e., tweets supposed to be read. Named entities are extracted via several third-party Natural Language Processing (NLP) Web services in parallel, and are then reported to Google Analytics, which is used to store, analyze, and compute trends by pivoting Analytics data, e.g., users’ geographic location, with the recorded named entities.

## 1 Introduction

The remainder of the paper is structured as follows:

### 1.1 Twitter Trends

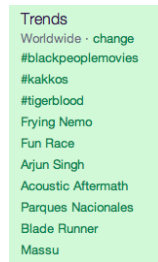
Twitter Trends was introduced by Twitter during the summer of 2008. This service was implemented to reflect “what’s happening in my world right now” on the micro-blogging platform. To compute trends, Twitter uses an algorithm which analyzes the words and hashtags in tweets. Although the concrete implementation details are kept secret, it is known<sup>2</sup> that the algorithm considers three main characteristics of a word/hashtag to determine if it is part of a “trend”:

---

<sup>1</sup> <http://blog.twitter.com/2008/09/twitter-trends-tip.html>

<sup>2</sup> <http://blog.twitter.com/2010/12/to-trend-or-not-to-trend.html>

- the quantity, i.e., the absolute number of appearances of a given word/hashtag among all users’ tweets.
- the velocity, i.e., the frequency of appearance of that word/hashtag. The higher the frequency, the more popular the word/hashtag.
- the newness, e.g., the (at the time of writing brand new) hashtag #ipad2 will be given priority over an old and generic tag like #awesome to be featured as “trend”. The freshness of trends is also assured by the fact that the frequency is computed over a recent set of tweets, although we have no further details of the exact computation period – Twitter keeps this information secret.



**Fig. 1.** Screenshot of Twitter Trends as of Friday, March 4, 2011, 4pm CET.

## 1.2 Google Chrome Extensions

Google Chrome extensions<sup>3</sup> are small software programs that can be installed to enrich the browsing experience with the Google Chrome browser. They are written using a combination of standard Web technologies, such as HTML, JavaScript, and CSS. Chrome extensions bundle all their resources into a single file that gets usually (but not necessarily) distributed through the Chrome Web Store. There are several types of extensions, for this paper we focus on extensions based on so-called content scripts. Content scripts are JavaScript programs that run in the context of Web pages, similar to the Firefox Greasemonkey extension<sup>4</sup>. By using the standard Document Object Model (DOM), they can read or modify details of the Web pages a user visits. Examples of such modifications are, e.g., changing hyperlinks to remove potential @target=”\_blank” attributes, or increasing the font size.

## 1.3 Google Analytics

Google Analytics is Google’s free Web analysis solution allowing for detailed statistics about the visitors of a website. The software is implemented by in-

<sup>3</sup> Google Chrome Extensions: <http://code.google.com/chrome/extensions/index.html>. Text partly adapted from the description to be found there.

<sup>4</sup> Firefox Greasemonkey extension: <http://www.greasemonkey.net/>

cluding the Google Analytics Tracking Code (GATC), an invisible snippet of JavaScript code that a webmaster needs to add onto the to-be-tracked pages of a website. This code collects visitor data by requesting a specific 1x1 pixel image on Google's servers, where the page and user data is encoded in the query part of the image's URL. In addition to that, the snippet sets a first party cookie on visitors' computers in order to store anonymous information such as the timestamp of the current visit, whether the visitor is a new or returning visitor, and the referrer of the website that the visitor came from. Part of the shared visitor information is the IP address, which allows for IP-based geolocation of visitors. Depending on the region, the accuracy of IP-based geolocation can be down to urban district level.

## 2 Twitter Swarm NLP Extension

With our Twitter Swarm NLP extension<sup>5</sup>, we inject JavaScript code via a content script into the Twitter.com homepage. By installing this extension, users explicitly opt-in to their data as a Twitter.com visitor being tracked by Google Analytics. The extension first checks if the user is logged in, and if so, retrieves the tweets of the logged-in user's timeline (<http://twitter.com/#>), or search result page (e.g., <http://twitter.com/#!/search/twitter>), or profile page (e.g., [http://twitter.com/#!/timberners\\_lee](http://twitter.com/#!/timberners_lee)) one-by-one, and performs NLP analysis via a remote NLP Web service (see section 2.1) on each of the tweets. The extracted entities are then displayed on the righthand-pane of the Twitter.com homepage (see the left part of figure 2), and sent to Google Analytics for further processing (see the right part of figure 2).

### 2.1 Twitter Swarm NLP Web Service

We have created a wrapper NLP Web service that merges results from existing third-party NLP Web services, namely from OpenCalais<sup>6</sup>, Zemanta<sup>7</sup>, AlchemyAPI<sup>8</sup>, and DBpedia Spotlight<sup>9</sup>. While the original calls to each particular NLP service are all HTTP POST-based, we have implemented the wrapper service GET- and POST-based. All NLP Web services return entities with their types and/or subtypes, names, relevance, and URIs that link into the LOD cloud. The problem is that each service has implemented its own type system, and providing mappings for all of them would be a rather time-consuming task. However, as all services offer links into the LOD cloud<sup>10</sup>, the desired type information can be pulled from there in a true Linked Data manner. The least common multiple of the results for the sample query "Google Translate" (i.e., the result of the call

<sup>5</sup> <https://chrome.google.com/webstore/detail/dpbphenfakflfmdlanimlemacankjol>

<sup>6</sup> <http://www.opencalais.com/>

<sup>7</sup> <http://www.opencalais.com/>

<sup>8</sup> <http://www.alchemyapi.com/>

<sup>9</sup> <http://dbpedia.org/spotlight>

<sup>10</sup> <http://lod-cloud.net/>

to the Web service at <http://tomayac.no.de/entity-extraction/combined/Google%20Translate>) is depicted below. For the sake of clarity, we just show one entity with two URIs while the original result contained seven entities among which six were directly relevant and one was closely related.

```
[
  {
    "name": "Google Translate",
    "relevance": 0.7128319999999999,
    "uris": [
      {
        "uri": "http://dbpedia.org/resource/Google_Translate",
        "source": "alchemyapi"
      },
      {
        "uri": "http://rdf.freebase.com/ns/en/google_translate",
        "source": "zemanta"
      }
    ],
    "source": "alchemyapi,zemanta"
  }
]
```

## 2.2 Technical Implementation Of the Twitter Swarm NLP Extension

Twitter.com is an Ajax-dependent website,<sup>11</sup> which makes use of so-called hashbang URIs. Currently the extension is implemented in a way to run once upon page load, however, not upon Ajax refreshes of the page. Each tweet gets sent one-by-one to the Twitter Swarm NLP Web service, which is described in section 2.1.

## 2.3 Dealing With Extracted Entities On the Client Side

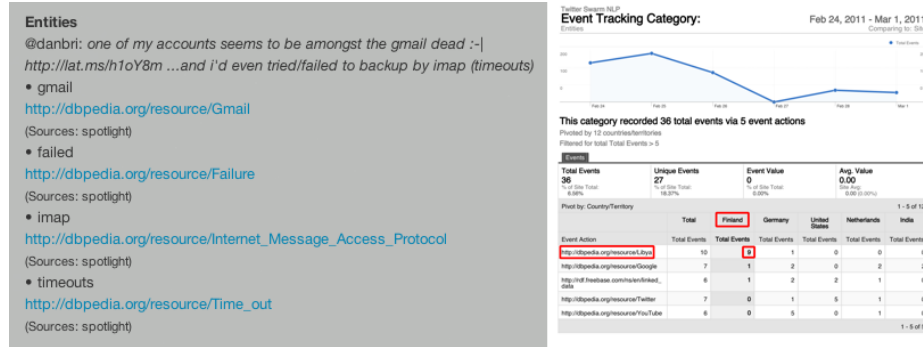
## 2.4 Dealing With Extracted Entities On the Google Analytics Side

# 3 Related Work

## 3.1 Linked Open Social Signals (TWARQL)

In the Linked Open Social Signals project [3] Mendes et al. investigate the representation of microposts as Linked Open Data (the authors call the opinions, observations, and suggestions contained in microposts “social signals”, hence the project name Linked Open Social Signals). Mendes et al. address the problem of information overload caused by the sheer amount of microposts. While the micropost community has come up with hashtags in order to categorize microposts, these hashtags are ambiguous and have to be explicitly added to the micropost by the author. Given the typical length limitations of microposts (often 140 characters), sometimes hashtags are left out in favor of more text. The project’s

<sup>11</sup> See <http://www.jenitennison.com/blog/node/154> for a detailed analysis of Twitter’s use of hashbang URIs



**Fig. 2.** Left: Screenshot of the extracted entites of a particular tweet as displayed by the Twitter Swarm NLP Extension. Right: Entities pivoted by country. The entity represented by the URL <http://dbpedia.org/resource/Lybia> appeared in nine tweets on timelines of users located in Finland (red borders in the screenshot).

main goal is thus to enable collective analysis of social signals for sensemaking by making use of Linked Open Data principles in combination with realtime push models. The approach consists of the following steps:<sup>12</sup>

- Extract content (entity mentions, hashtags and URLs) from microposts
- Encode content in RDF format using common vocabularies
- Enable SPARQL querying of microposts
- Enable subscription to micropost streams that match a given query
- Enable scalable real-time delivery of streaming data via SparqlPuSH

The authors maintain a client-side JavaScript application allowing for users to search for tweets matching a customizable SPARQL query or to subscribe to tweet streams in a realtime “push” way filtered according to the user’s request (so-called concept feeds).

### 3.2 Twitris 2.0

With Twitris 2.0 [2], Jadhav et al. present an application to find out what is when being said about an event, to detect how topics of discussion are changing over a period of time, and finally to check whether there are regional differences in the opinions on a given topic. The approach consists of picking trending hashtags from Twitter, which are then expanded by data from Google Insights for Search<sup>13</sup>. Using this set of search terms, find related hashtags by performing a Twitter search in order to detect topic drifts. In order to locate the origin of a tweet, Twitris uses the approximation of the location given in the user’s Twitter profile. This approach works well if there is geocodable data (like “Austin,

<sup>12</sup> Steps adapted from <http://wiki.knoesis.org/index.php/Twarql>

<sup>13</sup> <http://www.google.com/insights/search/>

Texas”), however, fails if there is generic data (like “somewhere under the rainbow”). Tweets about a given set of topics can then be examined on a map view, enriched by relevant photo and video content. For each tweet location hotspot a tag cloud with related tags is displayed, and the data can be sliced by time.

### 3.3 Semantic-MicroBlogging (SMOB)

SMOB [4] by Passant et al. is a Semantic MicroBlogging framework that enables a distributed, open and semantic microblogging experience based on Semantic Web and Linked Data technologies. Microposts get annotated with common vocabularies like FOAF<sup>14</sup> and SIOC<sup>15</sup>. SMOB relies on distributed hubs that communicate with each other to exchange microblog posts and subscriptions, which can also be cross-posted to Twitter. The authors suggest people to use meaningful hashtags such as `#dbp:Eiffel.Tower`, or `#geo:Paris.France`, in the style of RDF prefixes for DBpedia [1] and GeoNames. SMOB allows for manually annotating hashtags with URIs and RDF, with the objective of making microposts accessible for, e.g., lookup services such as Sindice [5].

### 3.4 Twopular

Twopular<sup>16</sup> is an experiment by Martin Dudek with the objective of analyzing current Twitter trends. Therefore, for a given set of at the particular moment current trends the most recent tweets are obtained, and in an interval of five minutes run through the OpenCalais Web service in order to find tags. By having tags for Twitter trends, another way of searching for trends – and more importantly the possibility to interrelate trends based on tag similarity – gets enabled. The author sees the feature more like a “linguistic experiment”<sup>17</sup>, however states that the first results seem to be promising. In our tests of the service we could affirm the author’s self-assessment, e.g., the Twitter trend (at the time of writing, March 7, 2011, 2PM CET) “Prince Andrew” was mapped to the OpenCalais tag “Prince Andrew, Duke of York”, and related tags were, among others, “Sex offender registration”, and “X-Offender”, where the story behind the trend was that people were tweeting about Prince Andrew of England, whose close friend was found to be a pedophile.

## 4 Conclusion

Contributions: time filters (via Analytics), geographical pivoting (via Analytics)

As seen in the Related Work section, semantic analysis of a (real-time) Twitter stream is not new and has been successfully exploited to analyse tweets produced by the Twitter community. What we propose here is an insight into tweets consumers’ interests to provide a more accurate view of Twitter trends.

<sup>14</sup> <http://www.foaf-project.org/>

<sup>15</sup> <http://sioc-project.org/>

<sup>16</sup> <http://twopular.com/>

<sup>17</sup> <http://twopular.com/blog/?p=308>



Fig. 3. Screenshot of the Twopular trends page.

## References

1. S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. G. Ives. DBpedia: A Nucleus for a Web of Open Data. In *6<sup>th</sup> International Semantic Web Conference (ISWC'07)*, pages 722–735, Busan, Korea, 2007.
2. A. Jadhav, H. Purohit, P. Kapanipathia, P. Ananthram, A. Ranabahu, V. Nguyen, P. N. Mendes, A. G. Smith, M. Cooney, and A. Sheth. Twitris 2.0 : Semantically empowered system for understanding perceptions from social data. Semantic Web Challenge at the *9<sup>th</sup> International Semantic Web Conference (ISWC2010)*, Shanghai, China, 2010.
3. P. N. Mendes, A. Passant, P. Kapanipathi, and A. P. Sheth. Linked open social signals. *Web Intelligence and Intelligent Agent Technology, IEEE/WIC/ACM International Conference on*, 1:224–231, 2010.
4. A. Passant, T. Hastrup, U. Bojars, and J. Breslin. Microblogging: A semantic and distributed approach. In *Proceedings of the 4<sup>th</sup> Workshop on Scripting for the Semantic Web, Tenerife, Spain, June 02, 2008, CEUR Workshop Proceedings*, 2008.
5. G. Tummarello, R. Delbru, and E. Oren. Sindice.com: weaving the open linked data. In *6<sup>th</sup> International Semantic Web Conference (ISWC'07)*, pages 552–565, Busan, Korea, 2007.