



**HAL**  
open science

## Characterizing Web Syndication Behavior and Content

Zeinab Hmedeh, Nelly Vouzoukidou, Nicolas Travers, Vassilis Christophides,  
Cédric Du Mouza, Michel Scholl

► **To cite this version:**

Zeinab Hmedeh, Nelly Vouzoukidou, Nicolas Travers, Vassilis Christophides, Cédric Du Mouza, et al.. Characterizing Web Syndication Behavior and Content. WISE'11, The 12th International Conference on Web Information System Engineering, Oct 2011, Sydney, Australia. pp.29-42. <hal-00737239>

**HAL Id: hal-00737239**

**<https://hal.science/hal-00737239v1>**

Submitted on 1 Oct 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# Characterizing Web Syndication Behavior and Content

Zeinab Hmedeh<sup>1</sup>, Nelly Vouzoukidou<sup>2</sup>, Nicolas Travers<sup>1</sup>,  
Vassilis Christophides<sup>2</sup>, Cedric du Mouza<sup>1</sup>, and Michel Scholl<sup>1</sup>

<sup>1</sup> CEDRIC Laboratory - CNAM - Paris, France  
firstname.lastname@cnam.fr

<sup>2</sup> FORTH/ICS and Univ. of Crete - Heraklion, Greece  
(christop,vuzukid)@ics.forth.gr

**Abstract.** We are witnessing a widespread of web syndication technologies such as RSS or Atom for a timely delivery of frequently updated Web content. Almost every personal weblog, news portal, or discussion forum employs nowadays RSS/Atom feeds for enhancing *pull-oriented* searching and browsing of web pages with *push-oriented* protocols of web content. Social media applications such as Twitter or Facebook also employ RSS for notifying users about the newly available posts of their preferred friends. Unfortunately, previous works on RSS/Atom statistical characteristics do not provide a precise and updated characterization of feeds' behavior and content, characterization which can be used to successfully benchmark effectiveness and efficiency of various RSS processing/analysis techniques. In this paper, we present the first thorough analysis of three complementary features of real-scale RSS feeds, namely, publication activity, items structure and length, as well as, vocabulary of its content which we believe are crucial for Web 2.0 applications.

**Keywords:** RSS/Atom Feeds, Publication activity, Items structure and length, textual vocabulary composition and evolution

## 1 Introduction

Web 2.0 technologies have transformed the Web from a publishing-only environment into a vibrant information place where yesterday's end users become nowadays content generators themselves. Web syndication formats such as RSS or Atom emerge as a popular mean for timely delivery of frequently updated Web content. According to these formats, information publishers provide brief summaries of the content they deliver on the Web, called *items*, while information consumers subscribe to a number of RSS/Atom *feeds* (i.e., streams or channels) and get informed about newly published items. Today, almost every personal weblog, news portal, or discussion forum employs RSS/Atom feeds for enhancing traditional *pull-oriented* searching and browsing of web pages with *push-oriented* protocols of web content. Furthermore, social media applications

such as Twitter or Facebook also employ RSS for notifying users about the newly available posts of their preferred friends (or followees).

Unfortunately, previous works on RSS/Atom statistical characteristics [15, 27, 13] do not provide a precise and updated characterization of feeds' behavior and content which could be effectively used for tuning refreshing policies of RSS aggregators [24, 22], benchmarking scalability and performance of RSS continuous monitoring mechanisms [19, 9, 7, 8, 25, 5] or comparing various techniques for RSS items mining, recommendation, enrichment and archiving [3, 26]. In this paper, we present the first thorough analysis of three complementary features of real-scale RSS/Atom feeds, namely, *publication activity*, *items structure and length*, as well as, *vocabulary of the textual content*.

Our empirical study relies on a large-scale testbed acquired over a 8 month campaign from March 2010 in the context of the French ANR project Roses<sup>3</sup>. We collected 10,794,285 items originating from 8,155 productive feeds (spanning over 2,930 different hosting sites) out of 12,611 feeds, without communication and validation errors against RSS/Atom specifications, harvested from major RSS/Atom directories, portals and search engines (such as syndic8.com, Google Reader, feedmil.com, completeRSS.com, etc.). Then, we have identified six representative types of information sources delivering their content as RSS/Atom feeds: **Press** (for newspapers and agencies), **Blogs** (for personal weblogs), **Forums** (for discussion and mailing lists), **Sales** (for marketing web sites), **Social media** (for micro-blogging such as twitter, digg, yahoo! groups, and blogosphere) and **Misc.** (*e.g.*, for news sites with medical or city/group information as well as podcasts). To the best of our knowledge, we are the first to detail the composition of our testbed whose type and number of sources is representative of the web syndication universe. The acquired RSS/Atom items are stored into a local warehouse<sup>4</sup> using MySQL, where for each feed only one item occurrence is kept. Since their *pubDate* was not made always available by the originating RSS/Atom feeds, we timestamped items based on their acquisition time in the warehouse. The main conclusions drawn from our experimental study are:

1. As analyzed in section 2, 17% of RSS/Atom feeds produce 97% of the items of the testbed. In their majority, productive feeds (*i.e.* with >10 items per day) exhibit a regular behavior without publication bursts, thus are more predictable in their publication behavior. As expected, micro-blogging feeds from social media are more productive than those from blogs while press sources lie in between. The average publication rate among all feeds of the testbed has been measured to be 3.59 items a day;
2. As highlighted in section 3, the most popular RSS/Atom textual elements are *title* and *description* while the average length of items is 52 terms (which has not been reported so far in related work). It is clearly advertisement greater than bids (4-5 terms [10]) or tweets (15 terms at most [21]) but smaller than blogs (250-300 terms [16]) or Web pages (450-500 terms excluding tags [14]).

---

<sup>3</sup> [www-bd.lip6.fr/roses](http://www-bd.lip6.fr/roses)

<sup>4</sup> Available on line at [deptmedia.cnam.fr/~traversn/roses](http://deptmedia.cnam.fr/~traversn/roses)

- In addition, re-publication of items across feeds is rare since we identified only 0.41% of duplicates among distinct feeds hosted by different sites;
- As studied in section 4, the total number of extracted terms from items written in English is 1,537,730 out of which only a small fraction (around 4%) is found in the WordNet dictionary. This is due to the heavy use in RSS/Atom textual elements of named entities (person and place names), URLs and email addresses as well as numerous typos or special-purpose jargon. We formally characterized the total vocabulary growth using Heaps' law, as well as the number of occurrences of the corresponding terms with a *stretched exponential distribution* used for the first time in the literature in this respect. We observed that the ranking of vocabulary terms does not significantly vary during the 8 month period of the study for frequent terms.

## 2 Feeds Analysis

In this Section we are interested in characterizing the composition of our testbed in terms of RSS/Atom feeds' type are originating from, as well as, in studying their publication activity. Although global statistics regarding Web2.0 activity are constantly monitored<sup>5</sup>, an in depth analysis of RSS feeds productivity was not already reported in the literature. Knowing that highly active feeds are more predictable in their publication behavior, would guide for instance resource allocation mechanisms to be tied to the feeds category.

### 2.1 Source Type

Six types of information sources delivering their content as RSS/Atom feeds were identified: **Press** (for newspapers and agencies), **Blogs** (for personal weblogs), **Forums** (for discussion and mailing lists), **Sales** (for marketing web sites), **Social media** (for micro-blogging such as twitter, digg, yahoo! groups, and blogosphere) and **Misc.** (*e.g.*, for news, medical, city/group information, podcasts). Then, the 8,155 productive feeds were classified under these six types using information available in feeds' title, category, link or in specific hosting sites.

Type	% of feeds	% of items	$\frac{\# \text{ items}}{\# \text{ feeds}}$
Social Media	1.77%	9.45%	7085.03
Press	9.99%	38.82%	5141.24
Forum	1.51%	3.62 %	3178.01
Sales	11.32%	15.49%	1811.92
Misc.	41.47%	25.47%	812.99
Blog	33.93%	7.14%	278.55

**Table 1.** Source types of RSS/Atom feeds

Table 1 depicts for each type, the corresponding percentage of feeds and items as well as the average number of items per feed. **Social media**, **Press**

<sup>5</sup> [thefuturebuzz.com/2009/01/12/social-media-web-20-internet-numbers-stats](http://thefuturebuzz.com/2009/01/12/social-media-web-20-internet-numbers-stats)

and **Forums** are more productive (with an average number of items per feed ranging from 3178.01 to 7085.03) than **Sales**, **Misc** and **Blogs** (with less than 2000 items on average per feed). As discussed in the following sections, feeds' behavior can be further refined by considering daily publication rates as well as the corresponding activity variability over time. We believe that the composition of our testbed is representative of the Web 2.0 universe. For instance the fact that blogs as opposed to Social media provide low productive feeds compared to their number, as well as the fact that numerous Press and Sales feeds are actually available, are reflected in the composition of our testbed.

## 2.2 Publication Activity

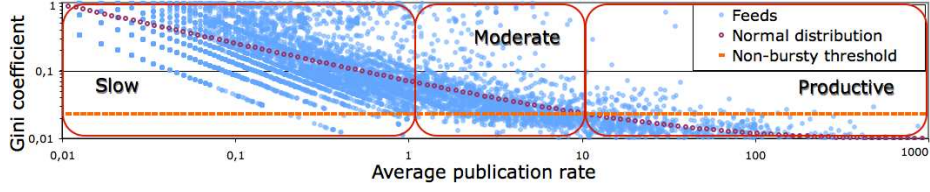
A time-agnostic characterization of RSS feeds activity has been originally proposed in [24] by studying the distribution of the number of feeds publishing at a given rate  $x$ . A similar power law behavior  $ax^b$  was observed for the 8,155 feeds, although with different coefficients  $a = 1.8 \times 10^3$  and  $b = -1.1$ . This is due to the presence in our testbed of more productive feeds than in [24] (featuring more blog-originating feeds). A coarse-grained characterization of the temporal variation of feeds activity has been suggested in [27] which analyzes publication burstiness. With a similar burst definition of at least 5 times the average publication rate during a unit of time (a day), 89% of our feeds produce bursts. However, the remaining 11% of feeds with bursts produce more than 81% of the items. Thus, this burstiness measure is not sufficient for a precise characterization of feeds activity.

We prefer the *Gini coefficient* [20], denoted as  $G$ , to characterize the variability of feeds' publication activity over time.  $G$  is adequate for time series (feeds are time series), since it does not only take into account the deviation from a mean value (as in the case of burstiness) but also the temporal variation of this deviation. It has been widely used for analyzing a variety of phenomena arising in economy, geography, engineering, or in computer science (e.g., self-join size, supervised learning).  $G$  is defined as:

$$G = \frac{\sum_{i=1}^n \sum_{j=1}^n |y_i - y_j|}{2 \times n \times \sum_{i=1}^n y_i}$$

with  $y_i$  denoting the number of items, sorted in increasing order, over the  $n$  days. A  $G$  value close to 1 suggests that the number of items of a feed significantly varies over time.

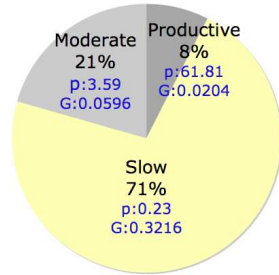
Figure 1 depicts in log-log-scale  $G$  vs the average publication rate for each of the 8,155 feeds. Feeds with a  $G$  value less than 0.02293 (below the horizontal dashed line) did not exhibit a single burst during the entire 8-month period. Three classes of feed activity are identified. The first one called "**productive class**", comprises 614 feeds which produce more than 10 items a day with a very low temporal variation ( $G$  less than 0.03). The second one called "**moderate class**", gathers 1,677 feeds that publish between 1 and 10 items per day with a



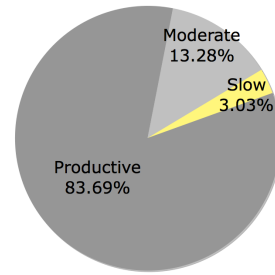
**Fig. 1.** Activity classes of RSS/Atom feeds

moderately low temporal variation ( $G$  less than 0.1). The third one called “slow class”, represents the large majority of the feeds. They publish less than 1 item a day (0.23 on average) and exhibit a strong temporal variation in the number of items ( $G = 0.32$  on average). The average  $G$  value as a function of the average publication rate  $p$  (dotted curve) is approximated by the following function:

$$G(p) = 5.53 \times 10^{-2} \times p^{-0.59} + 4.48 \times 10^{-3} \times p^{0.101}$$



**Fig. 2.** Feeds per activity class



**Fig. 3.** Items per activity class

Unlike [15] reporting that 57% of the feeds have a publication rate  $> 24$  items a day, in Figure 2 we can see that only 8% of our 8,155 feeds were actually very productive (with  $> 10$  items a day) while 71% had a low publication activity. Note also that from the initially harvested 12,611 feeds around 35% did not publish any item during the entire 8-month period. Of course productive feeds, even though they represent a minority, account for the majority of the items harvested. We see in Figure 3 that feeds of **productive** class publish 83.69% of the total number of items while the majority of feeds from **slow** class produce only 3.03% of the items (Power Law behavior of feeds’ publication rate).

Type	productive class			moderate class			slow class		
	%	p	G	%	p	G	%	p	G
Social Media	44.9%	65.33	.015	8.1%	5.57	.026	47.0%	0.18	.470
Press	27.4%	70.76	.020	43.0%	3.61	.036	29.6%	0.34	.237
Forum	21.0%	54.51	.018	20.2%	3.85	.042	58.8%	0.29	.348
Sales	8.3%	83.64	.022	19.0%	3.43	.122	72.7%	0.23	.454
Misc	2.7%	63.12	.024	12.9%	3.70	.059	84.4%	0.23	.338
Blog	4.3%	15.53	.019	24.7%	3.37	.056	71%	0.22	.267

**Table 2.** Source types and activity classes of feeds

Last but not least, Table 2 provides for each source type and activity class, the corresponding percentage of feeds, the average publication rate ( $p$ ), as well

as, the average Gini coefficient (G). Feeds from Social Media are almost evenly distributed in **productive** and **slow** classes with a temporal variation in their publication rate which is more important in the latter than in the former class. This is related to *Twitter* like behavior with *Followers* and *Followees* (users can follow - slow class - or users can be followed and produce a lot - productive class). Press feeds exhibit a **moderate** and almost *regular* publication activity, while feeds originating from Forums, Sales, Misc and Blogs sites mainly exhibit a **slow** publication activity. It is worth also noticing that in this class, Sales are mostly bursty feeds while Blogs and Forums exhibit a more regular behavior.

### 3 Items Analysis

This section successively focuses on the analysis of items' structure and length as well as on the items' replication rate across RSS/Atom feeds which could be exploited for benchmarking scalability and performance of RSS continuous monitoring mechanisms.

#### 3.1 Item Structure

The empirical analysis presented in Table 3, reveals that a great number of fields (tags) foreseen in the XML specification of RSS or Atom formats are actually not utilized in items. While `title`, `description` and `link` are present in almost all the items, `pubDate` is missing in around 20% of the items, and the language information is missing in 30% of the feeds. Almost 2/3 of the items are not *categorized* while `author` information is present in less than 8% of the items. It is worth noticing that 16% of fields contain errors or are used with a wrong format (mixing RSS or Atom fields). Other XML tags are only sparsely used in our testbed and are not reported here. In a nutshell, RSS/Atom items are characterized by the predominance of textual information as provided by `title` and `description` over factual information provided by `category` and `author`.

<b>title</b>	<b>link</b>	<b>pubDate</b>	<b>desc.</b>	<b>Language</b>
99.82%	99.88%	80.01%	98.09%	69.14% (feed)
<b>author</b>	<b>category</b>	<b>GUID</b>	<b>ext.</b>	<b>RSS/Atom</b>
7.51%	33.94%	69.50%	29.73%	16.48%

**Table 3.** Popularity of XML tags in RSS items

#### 3.2 Item Length

Next, we focus on the length of items measured as the number of terms of the textual `title` and `description` fields. Items are short with a 52-terms size on the average (see Table 4). The high variance (11,885) of items length is mainly due to the large diversity of the `description` fields that can be either missing (or be a simple *url*) or oppositely be a long text (even an entire HTML document).

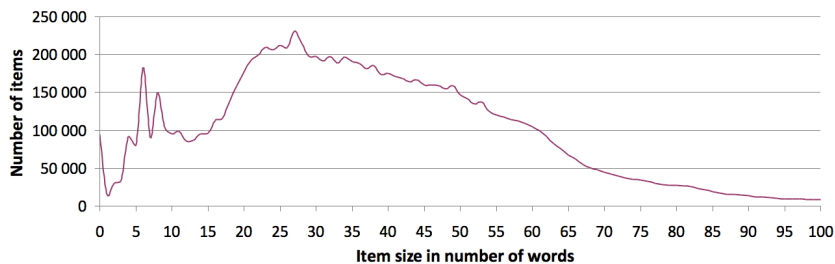
Figure 4 plots the number of items vs their length. A long-tail curve is observed in items length distribution as also reported in the literature for the size

	title + description	title	description
<b>Average</b>	52.37	6.81	45.56
<b>Max</b>	10,262	235	10,256
<b>Variance</b>	11,885.90	12.97	11,821.36

**Table 4.** Textual content characteristics in items

of Web documents [28]. 51.39% of the items have a length between 21 and 50 terms, and 14% between 8 and 20 terms. The peaks for length 6 and 8 are mainly due to **Sales** feeds (producing >55% of the items for these lengths) whose items respect a fixed-length pattern (*i.e.* items differ only by one or two terms).

The main conclusion from this analysis is that RSS/Atom items are longer than average advertisement bids (4-5 terms [10]) or tweets ( $\sim 15$  terms at most [21]) but smaller than the original blog posts (250-300 terms [16]) or Web pages (450-500 terms excluding tags [14]).



**Fig. 4.** Number of items per item length

### 3.3 Item Replication

Finally, we stick our attention to the replication of items either by feeds hosted by the same site (*intra-host* replication) or across different feeds (*inter-feed* replication). Replication detection is performed by exact matching of the item content based on a hash function. Out of the originally harvested 27 millions of items, 10.7 million distinct items per feed were extracted. Eliminating *intra-host* and *inter-feed* replication makes this number drop to 9,771,524 of items.

% of replication	< 10%	10 - 19%	20 - 29%	$\geq 30\%$
% of hosts	95.19%	1.09%	0.68%	3.04%
% of items	71.31%	10.88%	10.23%	7.58%

**Table 5.** % Hosts/items per intra-host replication

Table 5 reports the importance of *intra-host* replication that was measured for the 2,930 hosts (with distinct IPs) of the 8,155 feeds. 95% of the hosts (which publish 70% of the items) have less than 10% of replication. Only a few hosts (less than 5%) of **Sales**, **Press** (especially news agencies) and **Blogs** feeds publish most of the replicated items (replication rate  $> 30\%$ ).

Once intra-host replicates had been removed, *inter-feed* item replication was measured. Table 6 reports for each replicated item the number of *distinct* feeds in which it appears. Clearly replication across feeds in different hosts is negligible: it accounts for less than 0.5% of the total number of items. This behavior

# distinct feeds	1	2	$\geq 3$
% items	99.51%	0.41%	0.08%

**Table 6.** % of items replicated across feeds

can be explained by the absence in the testbed of RSS aggregators such as GoogleReader<sup>6</sup> or Yahoo Pipes<sup>7</sup> which systematically replicate other feeds.

## 4 Vocabulary Analysis

This Section focuses on the analysis of the vocabulary of terms extracted from the title and description fields of RSS/Atom items. In order to deduce information regarding the quality of the employed vocabulary (valid terms/typos), we restrict our analysis to the English language, since a large majority of the analysed items in our testbed was written in English. In addition, as previous studies report, a similar behavior to that of the English language is exhibited for text corpora written in different languages (e.g. French, Korean [4], Greek [6], etc.). An automatic filtering of feeds (items) written in English turned out to be a difficult task given that language information was not always available.

7,691,008 “English” items were extracted as follows: an item is considered as English if the language tag exists in the feed metadata and is set to “en”, or the feed url belongs to an English spoken host (*e.g.* “us” or “uk”), or a “.com” feed without any language tag. Then, a lexical analysis of items’ textual contents was conducted using standard stemming tools (such as SnowBall<sup>8</sup>) and dictionaries (such as WordNet<sup>9</sup>) for the English language. In particular, we distinguish between  $V_W$ , the vocabulary of (61,845) terms appearing in the WordNet dictionary, from  $V_{\overline{W}}$ , the vocabulary of remaining (1,475,885) terms composed mostly of jargon, named entities and typos. Stop-words<sup>10</sup>, *urls* and *e-mail* addresses were excluded and for  $V_{\overline{W}}$  the stemmed version of terms was kept while for  $V_W$  the terms in their base form as given by WordNet was added.

We believe that such a detailed characterization of feeds content will be beneficial to several Web 2.0 applications. For instance, knowledge regarding terms (co-)occurrence distribution could be helpful for efficient compression or indexing techniques of items textual content. The knowledge that the employed vocabulary is essentially composed of misspellings, spoken acronyms and morphological variants will greatly affect the choice of effective ranking functions for filtering the incoming items while it will enable us to devise realistic workloads for measuring the scalability and performance of Web 2.0 analysis systems (e.g., keyword tracking, buzz measuring, keyword-association in online consumer intelligence).

<sup>6</sup> [www.google.com/reader](http://www.google.com/reader)

<sup>7</sup> [pipes.yahoo.com](http://pipes.yahoo.com)

<sup>8</sup> [snowball.tartarus.org](http://snowball.tartarus.org)

<sup>9</sup> [wordnet.princeton.edu](http://wordnet.princeton.edu)

<sup>10</sup> [www.lextek.com/manuals/onix/stopwords2.html](http://www.lextek.com/manuals/onix/stopwords2.html)

#### 4.1 Term Occurrences

The global vocabulary ( $V=V_{\mathcal{W}}\cup V_{\overline{\mathcal{W}}}$ ) extracted from the English items reaches 1,537,730 terms. Figure 5 depicts the occurrences of terms belonging to  $V_{\mathcal{W}}$  and to  $V_{\overline{\mathcal{W}}}$  in decreasing order of their rank (frequency) in their respective vocabulary. As expected, terms from  $V_{\mathcal{W}}$  are much more frequent than terms from  $V_{\overline{\mathcal{W}}}$ . The multiple occurrences in the items of the 5,000 most frequent terms of  $V_{\mathcal{W}}$  (around 8% of its size) represent 87% of the total number of term occurrences from  $V_{\mathcal{W}}$ . The percentage of  $V_{\mathcal{W}}$  terms appearing in the most frequent terms of  $V$  drops quickly: from 90% in the 5,000 first terms, to 78% for the next 5,000, to 50% after 20,000 terms and to only 3% for the remaining 1.5 million terms.

In the following, we are interested in characterizing formally the distribution of  $V_{\mathcal{W}}$  and  $V_{\overline{\mathcal{W}}}$  terms' occurrences. In this respect, Zipf's law distributions have been traditionally used in the literature [2, 17] for various text corpora:

$$f(r) = \frac{K}{r^\theta}$$

where  $r$  is the term rank and  $\theta$  and  $K$  are constants.

However, as can be seen in Figure 5 the corresponding curve for  $V_{\mathcal{W}}$  has a significant deviation from Zipf's law, *i.e.*, from a straight line in log-log scale. This deviation is smaller for the  $V_{\overline{\mathcal{W}}}$  curve. Similar deviations have been already reported for web related text collections [2, 17, 1, 28, 10] and few attempts have been made to devise more adequate distributions. [18] tried to generalize the Zipf's law by proposing a Zipf - Mandelbrot distribution while [13] suggested to use a Modified Power Law distribution. Although the latter laws can approximate the slow decrease at the beginning of the curves, their tail in log-log scale is close to a straight line, which does not fit neither to our terms' occurrences nor to any of the distributions in the aforementioned studies.

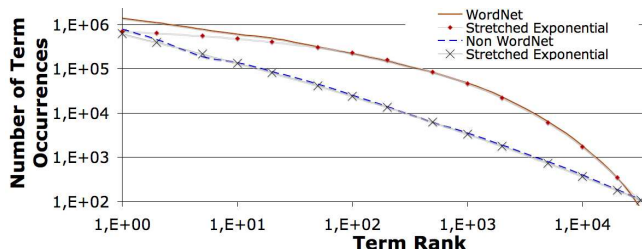


Fig. 5. Occurrences of terms from  $V_{\mathcal{W}}$  and  $V_{\overline{\mathcal{W}}}$

A *stretched exponential* distribution was found to better capture the tail of the distributions for both vocabularies. It is defined as follows:

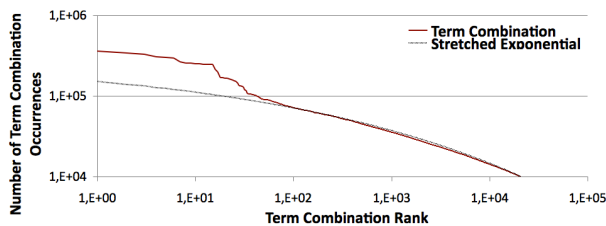
$$f(r) = K \times e^{-(r/c)^\theta}$$

where  $\theta$  expresses how fast the slope of the curve decreases, while constant  $c$  defines the curvature: for values closer to zero the distribution is closer to a straight line in log-log scale. To the best of our knowledge the stretched exponential distribution has been so far utilized to study physics phenomena but has

never been associated with term distributions before, even though [12] suggests this distribution as an alternative to power laws. The constants that lead to the best fitting curves were obtained using the chi-square test (see Table 7). Lower values of chi-square indicate a better fit of the distribution to the empirical data.

	Zipf	Zipf Mandelbrot	Modified power law	Stretched Exponential
$V_W$	2,845	2,493	966	<b>49.8</b>
$\overline{V_W}$	11.5	11.5	9.2	<b>9</b>
Term co-occurrence	2,353	1,543	31,653	<b>45</b>

**Table 7.** Chi-square for 4 distributions (deg. of freedom: 210 for  $V_W$  and 200 for  $\overline{V_W}$ )



**Fig. 6.** Co-occurrences of terms from  $V$

Figure 6 represents the occurrences of the 32,000 most frequent term combinations in items. These combinations are all pairs of terms, not necessarily consecutive, that appear in the item’s title or description regardless of the vocabulary they belong to. In contrast to [10], where the same test was performed for sponsored data, we find that this distribution does not follow the Zipf’s law. In Figure 6 the curve has been fitted using the stretched exponential distribution whose chi-square value is the lowest one (see table 7). The divergence of the empirical data from the stretched exponential distribution given above is limited to the 32 most frequent combinations of terms, which appear almost 4 times more than expected. These combinations are essentially due to high interests for some serials (*e.g.* like **Lost-fan**), a behavior also observed in [10].

Occurrences	1	2	3	4	$\geq 5$
# terms	659,159	261,906	124,450	77,814	414,401

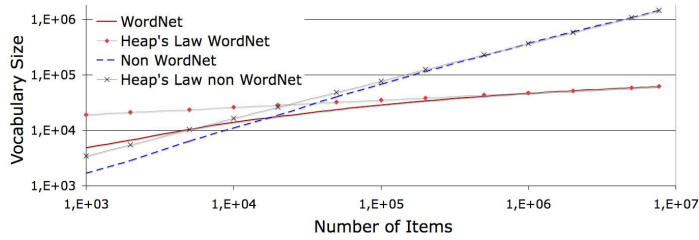
**Table 8.** Terms’ number per occurrences

As the number of occurrences varies widely from one term to another (see Figure 5) in Table 8 we provide a summary of term occurrences distribution. Terms with less than 5 occurrences represent roughly 76% of the global vocabulary size while terms with only one occurrence account for 45% of the vocabulary. A qualitative analysis over two samples of 5,000 terms comprising respectively the most and the less frequent terms of  $\overline{V_W}$  showed that 61% of most frequent terms are named-entities and 27% are acronyms or shortcuts. On the other hand, less frequent terms are named-entities (35%), foreign-words (22%), composed and concatenation of words (22%) and misspelling (14%). Clearly, language of

non-authoritative web document collections on diverse topics is subject to many kinds of imperfections and errors [1], like typos, mistake or slang, but also special-purpose jargon (*e.g.* technical terms in medicine or computer science).

## 4.2 Vocabulary size

Figure 7 illustrates the size growth wrt the number of items progressively stored in our warehouse for the two vocabularies of terms. Clearly, the size of  $V_{\overline{\mathcal{W}}}$  evolves much faster than that of  $V_{\mathcal{W}}$ : at the end of the 8-month period, it is 24 times bigger than the size of the former ( $|V_{\mathcal{W}}| = 61,845$  while  $|V_{\overline{\mathcal{W}}}| = 1,475,885$ ).



**Fig. 7.** Evolution of the two vocabularies

For  $V_{\mathcal{W}}$ , a rapid increase is observed for the first 140,000 items to reach the 31,000 most frequent terms of  $V_{\mathcal{W}}$ , while the size barely doubles up to the end, with the inclusion of less frequent terms. On the other hand,  $V_{\overline{\mathcal{W}}}$  has a smoother growth: for 140,000 items  $V_{\overline{\mathcal{W}}}$  comprises 88,000 terms, *i.e.* only 6% of the total  $V_{\overline{\mathcal{W}}}$  size which reaches almost 1.5 million terms after processing all items.

The vocabulary size growth is usually characterized by a Heap's law distribution [2, 17] (see Figure 7), defined as follows:

$$|V(n)| = K \times n^\beta$$

where  $n$  is the number of collected items while  $K$  and  $\beta$  (values in  $[0, 1]$ ) are constants depending strongly on the characteristics of the analyzed text corpora and on the model used to extract the terms [28].  $\beta$  determines how fast the vocabulary evolves over time with typical values ranging between 0.4 and 0.6 for medium-size and homogeneous text collections [2]. [28] reports a Heap's law exponent lying outside this range ( $\beta=0.16$ ) for a 500 MB collection of documents from the Wall Street Journal. Table 4.2 specifies the constants chosen for Heap's laws approximating the global vocabulary growth as well as  $V_{\mathcal{W}}$  and  $V_{\overline{\mathcal{W}}}$ .

Clearly, the Heap's law exponent ( $\beta$ ) of the global vocabulary is affected by the evolution of  $V_{\overline{\mathcal{W}}}$  rather than by  $V_{\mathcal{W}}$  whose size is significantly smaller (attributed to the slow acquisition of less commonly used terms). The exponent for  $V_{\overline{\mathcal{W}}}$  (0.675) slightly higher than those reported in the literature [2, 17] indicates a faster increase of the vocabulary size due the aforementioned language imperfections of the items in our testbed. It should be stressed that this behavior is also exhibited by the evolution of vocabularies related to other user-generated

	$ V_w + V_{\overline{w}} $	$ V_w $	$ V_{\overline{w}} $
$K$	51	7,921	33
$\beta$	0.65	0.13	0.675

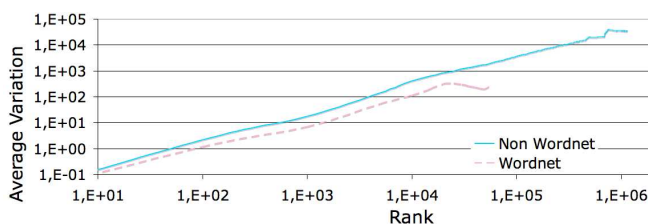
**Table 9.** Heap Laws constants

textual content (such as web queries [29]). Finally, the high Heap’s law coefficient ( $K$ ) for  $V_w$  is explained by the rapid vocabulary growth in the beginning due to the fast acquisition of the very popular terms.

This behavior is also observed in web queries vocabulary evolution [29, 23] where the exponents of the Heap’s law exceed the upper bound of 0.6. In the matching Heap’s law of [29], the exponent of the vocabulary distribution for web queries was 0.69, while in [23] it was 0.8136. In both items and web queries, the text is sent once by the users or publishers and is never corrected (i.e. in a case of misspelling), which leads to a higher number of errors and therefore a faster than expected vocabulary evolution. Web documents, on the other hand, are frequently updated and possible misspellings are corrected.

### 4.3 Rank variation

Since term occurrences evolves over time as new items arrive, their corresponding ranking could be also subject to change. Figure 8 illustrates the average variation of ranks during the last month of our study. A rank variation is measured by the difference of a term rank between week  $t-1$  and week  $t$ . Obviously, this variation is proportional to the rank, *i.e.* less frequent terms are more likely to change their rank compared to the vocabulary of the previous week.



**Fig. 8.** Cumulative rank distance in  $V$  (last week)

We are finally interested in studying the weekly rank variation for specific rank ranges of  $V_{\overline{w}}$  terms for the 8 months of our study using the *Spearman* metric [11]. The *Spearman* metric is the sum of rank variations for a range of ranks from a week to another:

$$S(t) = \sum_{i=1}^n |r(t)_i - r(t-1)_i|$$

where  $t$  is a given week,  $r(t)_i$  and  $r(t-1)_i$  are the rank of the term  $i$  at week  $t$  and  $t-1$  respectively. Figure 9 depicts the weekly *Spearman* value for three classes of terms’ ranks: (a) highly-frequent terms (ranks between 1 and 250), (b) commonly-used terms (10,000 to 10,250); and (c) rare ones (500,000



**Fig. 9.** Weekly Spearman Variation for three rank ranges of  $V_{\overline{w}}$  terms

to 500,250). Note that the reference vocabulary for this experiment is chosen as the vocabulary obtained after the first three months of items’ acquisition.

Surprisingly enough, we have observed that the ranking of vocabulary terms does not significantly vary during the 8-month period of the study for highly-frequent terms. As a matter of fact, the sum of rank variations of highly-frequent terms is 4 orders of magnitude less than the corresponding sum for rare terms (where commonly-used terms lie in between). We notice that all Spearman values stabilize after 20 weeks. This can be justified by the fact that the number of occurrences of each term decreases the impact of incoming terms on the rank.

## 5 Summary

In this paper we presented an in-depth analysis of a large testbed of 8,155 RSS feeds comprising 10.7 million of items collected over a 8-month period campaign. We proposed a characterization of feeds according to their publication rate and temporal variability highlighting three different activity classes. In addition, we focused on RSS items length, structure and replication rate. Last but not least, a detailed study of the vocabulary employed in a significant subset of the RSS items written in English was conducted along three lines: occurrence distribution of the terms, evolution of the vocabulary size and evolution of term ranks.

## References

1. AHMAD, F., AND KONDRAK, G. Learning a Spelling Error Model from Search Query Logs. In *EMNLP* (2005).
2. BAEZA-YATES, R. A., AND RIBEIRO-NETO, B. A. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
3. BOURAS, C., POULOPOULOS, V., AND TSOGKAS, V. Creating Dynamic, Personalized RSS Summaries. In *ICDM* (2008), pp. 1–15.
4. CHOI, S.-W. Some statistical properties and zipf’s law in korean text corpus. *JQL* 7, 1 (2000), 19–30.
5. HAGHANI, P., MICHEL, S., AND ABERER, K. The gist of everything new: personalized top-k processing over web 2.0 streams. In *CIKM* (2010), ACM, pp. 489–498.
6. HATZIGEORGIU, N., MIKROS, G., AND CARAYANNIS, G. Word length, word frequencies and zipf’s law in the greek language. *JQL* 8, 3 (2001), 175–185.
7. HRISTIDIS, V., VALDIVIA, O., VLACHOS, M., AND YU, P. S. A System for Keyword Search on Textual Streams. In *SDM* (2007).

8. HU, C.-L., AND CHOU, C.-K. RSS Watchdog: an Instant Event Monitor on Real Online News Streams. In *CIKM* (2009), pp. 2097–2098.
9. IRMAK, U., MIHAYLOV, S., SUEL, T., GANGULY, S., AND IZMAILOV, R. Efficient Query Subscription Processing for Prospective Search Engines. In *USENIX* (2006), pp. 375–380.
10. KÖNIG, A. C., CHURCH, K. W., AND MARKOV, M. A Data Structure for Sponsored Search. In *ICDE* (2009), pp. 90–101.
11. KUMAR, R., AND VASSILVITSKII, S. Generalized distances between rankings. In *WWW* (2010), pp. 571–580.
12. LAHERRÈRE, J., AND SORNETTE, D. Stretched exponential distributions in nature and economy: "fat tails" with characteristic scales. *Eur. Phys. J. B* 2, 4 (1998), 525–539.
13. LAMBIOTTE, R., AUSLOOS, M., AND THELWALL, M. Word Statistics in Blogs and RSS Feeds: Towards Empirical Universal Evidence. *CoRR* (2007).
14. LEVERING, R., AND CUTLER, M. The portrait of a common html web page. In *ACM Symp. on Document Engineering* (2006), pp. 198–204.
15. LIU, H., RAMASUBRAMANIAN, V., AND SIRER, E. G. Client Behavior and Feed Characteristics of RSS, a Publish-Subscribe System for Web Micronews. In *IMC* (2005), pp. 3–3.
16. MA, S., AND ZHANG, Q. A Study on Content and Management Style of Corporate Blogs. In *HCI (15)* (2007), pp. 116–123.
17. MANNING, C. D., RAGHAVAN, P., AND SCHÜTZE, H. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
18. MONTEMURRO, M. A. Beyond the zipf-mandelbrot law in quantitative linguistics. *Physica A* 300, 3-4 (2001), 567–578.
19. PETROVIC, M., LIU, H., AND JACOBSEN, H.-A. CMS-ToPSS: Efficient Dissemination of RSS Documents. In *VLDB* (2005), pp. 1279–1282.
20. PITOURA, T., AND TRIANTAFILLOU, P. Self-join size estimation in large-scale distributed data systems. In *ICDE* (2008), pp. 764–773.
21. PRESS, O. U. Rt this: Oup dictionary team monitors twitterer's tweets, June 2009.
22. ROITMAN, H., CARMEL, D., AND YOM-TOV, E. Maintaining dynamic channel profiles on the web. *VLDB* 1, 1 (2008), 151–162.
23. SCHMIDT-MAENZ, N., AND KOCH, M. Patterns in search queries. In *Data Analysis and Decision Support* (2005).
24. SIA, K. C., CHO, J., AND CHO, H.-K. Efficient monitoring algorithm for fast news alerts. *TKDE* 19 (2007), 950–961.
25. SILBERSTEIN, A., TERRACE, J., COOPER, B. F., AND RAMAKRISHNAN, R. Feeding frenzy: selectively materializing users' event feeds. In *SIGMOD* (2010), pp. 831–842.
26. TADDESSE, F. G., TEKLI, J., CHBEIR, R., VIVIANI, M., AND YETONGNON, K. Semantic-Based Merging of RSS Items. *WWW* 13, 1-2 (2010), 169–207.
27. THELWALL, M., PRABOWO, R., AND FAIRCLOUGH, R. Are Raw RSS Feeds Suitable for Broad Issue Scanning? A Science Concern Case Study. *JASIST* 57, 12 (2006), 1644–1654.
28. WILLIAMS, H. E., AND ZOBEL, J. Searchable words on the Web. *JODL* 5, 2 (2005), 99–105.
29. ZIEN, J. Y., MEYER, J., TOMLIN, J. A., AND LIU, J. Web Query Characteristics and their Implications on Search Engines. In *WWW* (2001).