

## From the Cover: Semiotic dynamics and collaborative tagging

Ciro Cattuto, Vittorio Loreto, and Luciano Pietronero

*PNAS* 2007;104;1461-1464; originally published online Jan 23, 2007;  
doi:10.1073/pnas.0610487104

**This information is current as of May 2007.**

<b>Online Information &amp; Services</b>	High-resolution figures, a citation map, links to PubMed and Google Scholar, etc., can be found at: <a href="http://www.pnas.org/cgi/content/full/104/5/1461">www.pnas.org/cgi/content/full/104/5/1461</a>
<b>Related Articles</b>	A related article has been published: <a href="http://www.pnas.org/cgi/content/full/104/5/1443">www.pnas.org/cgi/content/full/104/5/1443</a>
<b>Supplementary Material</b>	Supplementary material can be found at: <a href="http://www.pnas.org/cgi/content/full/0610487104/DC1">www.pnas.org/cgi/content/full/0610487104/DC1</a>
<b>References</b>	This article cites 13 articles, 1 of which you can access for free at: <a href="http://www.pnas.org/cgi/content/full/104/5/1461#BIBL">www.pnas.org/cgi/content/full/104/5/1461#BIBL</a>  This article has been cited by other articles: <a href="http://www.pnas.org/cgi/content/full/104/5/1461#otherarticles">www.pnas.org/cgi/content/full/104/5/1461#otherarticles</a>
<b>E-mail Alerts</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .
<b>Rights &amp; Permissions</b>	To reproduce this article in part (figures, tables) or in entirety, see: <a href="http://www.pnas.org/misc/rightperm.shtml">www.pnas.org/misc/rightperm.shtml</a>
<b>Reprints</b>	To order reprints, see: <a href="http://www.pnas.org/misc/reprints.shtml">www.pnas.org/misc/reprints.shtml</a>

Notes:

# Semiotic dynamics and collaborative tagging

Ciro Cattuto<sup>\*†</sup>, Vittorio Loreto<sup>†‡</sup>, and Luciano Pietronero<sup>†</sup>

<sup>\*</sup>Museo Storico della Fisica e Centro Studi e Ricerche Enrico Fermi, Compendio Viminale, 00184 Rome, Italy; and <sup>†</sup>Dipartimento di Fisica, Università "La Sapienza," Piazzale Aldo Moro 2, 00185 Rome, Italy

Communicated by Nicola Cabibbo, University of Rome, Rome, Italy, November 30, 2006 (received for review June 20, 2006)

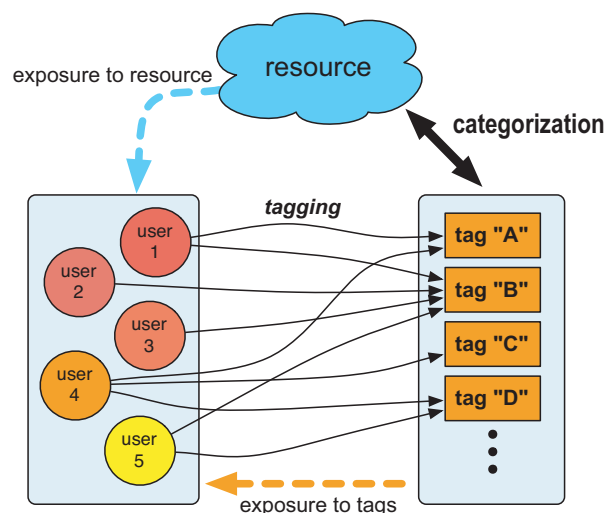
Collaborative tagging has been quickly gaining ground because of its ability to recruit the activity of web users into effectively organizing and sharing vast amounts of information. Here we collect data from a popular system and investigate the statistical properties of tag cooccurrence. We introduce a stochastic model of user behavior embodying two main aspects of collaborative tagging: (i) a frequency-bias mechanism related to the idea that users are exposed to each other's tagging activity; (ii) a notion of memory, or aging of resources, in the form of a heavy-tailed access to the past state of the system. Remarkably, our simple modeling is able to account quantitatively for the observed experimental features with a surprisingly high accuracy. This points in the direction of a universal behavior of users who, despite the complexity of their own cognitive processes and the uncoordinated and selfish nature of their tagging activity, appear to follow simple activity patterns.

online social communities | statistical physics | social bookmarking | information dynamics

Recently, a new paradigm has been quickly gaining ground on the World Wide Web: collaborative tagging (1–3). In web applications like *Del.icio.us* (<http://del.icio.us>), *Flickr* ([www.flickr.com](http://www.flickr.com)), *CiteULike* ([www.citeulike.org](http://www.citeulike.org)), and *Connotea* ([www.connotea.org](http://www.connotea.org)), users manage, share, and browse collections of online resources by enriching them with semantically meaningful information in the form of freely chosen text labels (tags). The paradigm of collaborative tagging has been successfully deployed in web applications designed to organize and share diverse online resources such as bookmarks, digital photographs, academic papers, music, and more. Web users interact with a collaborative tagging system by posting content (resources) into the system, and associating text strings (tags) with that content, as shown in Fig. 1. At the global level, the set of tags, although determined with no explicit coordination, evolves in time and leads toward patterns of terminology usage that are shared by the entire user community. Hence, one observes the emergence of a loose categorization system that can be effectively used to navigate through a large and heterogeneous body of resources.

Focusing on tags as basic dynamical entities, the process of collaborative tagging falls within the scope of semiotic dynamics (4–6), a new field that studies how populations of humans or agents can establish and share semiotic systems, typically driven by their use in communication. Indeed, the emergence of a folksonomy exhibits dynamical aspects also observed in human languages (7, 8), such as the crystallization of naming conventions, competition between terms, takeovers by neologisms, and more.

In the following, we adopt the point of view of complex systems science and try to understand how the “microscopic” tagging activity of users causes the emergence of the high-level features we observe for the ensuing folksonomy. We ground our analysis on actual tagging data extracted from *Del.icio.us* and *Connotea* and use standard statistical tools to gain insights into the underlying tagging dynamics. Based on this, we introduce a simple stochastic model for the tagging behavior of an “average” user, and show that such a model, despite its simplicity, is able to reproduce extremely well some of the observed properties. We



**Fig. 1.** Collaborative tagging. Schematic depiction of the collaborative tagging process: web users are exposed to a resource and freely associate tags with it. Their interaction with the system also exposes them to tags previously entered by themselves and by other users. The aggregated activity of users leads to an emergent categorization of resources in terms of tags shared by a community.

close giving an interpretation of the model parameters and pointing out directions for future research.

## Results

**Data Analysis.** The activity of users interacting with a collaborative tagging system consists of either navigating the existing body of resources by using tags, or adding new resources to the system. To add a new resource to the system, the user is prompted for a reference to the resource and a set of tags to associate with it. Thus the basic unit of information in a collaborative tagging system is a (user, resource, {tags}) triple, here referred to as post. Tagging events build a tripartite graph (with partitions corresponding to users, resources and tags, respectively) and such a graph, commonly referred to as folksonomy, can be subsequently used as a navigation aid in browsing tagged information (see Fig. 2). Usually, a post contains also a temporal marker indicating the physical time of the tagging event, so that temporal ordering can be preserved in storing and retrieving posts. Here we analyze data from *Del.icio.us* and *Connotea* and investigate the statistical properties of tag association. Specifically, we select a semantic context by extracting the resources associated with a given tag  $X$  and study the statistical distribution of tags cooccurring with  $X$

Author contributions: C.C., V.L., and L.P. designed research; C.C. and V.L. performed research; C.C. and V.L. contributed new reagents/analytic tools; C.C. analyzed data; and C.C., V.L., and L.P. wrote the paper.

The authors declare no conflict of interest.

<sup>†</sup>To whom correspondence should be addressed. E-mail: [vittorio.loreto@roma1.infn.it](mailto:vittorio.loreto@roma1.infn.it).

This article contains supporting information online at [www.pnas.org/cgi/content/full/0610487104/DC1](http://www.pnas.org/cgi/content/full/0610487104/DC1).

© 2007 by The National Academy of Sciences of the USA

ajax apple art article blog blogging blogs books browser business code community  
 computer cool css culture database design development diy download firefox  
 flash flickr framework free freeware fun gallery games google graphics hacks howto  
 html humor illustration images imported inspiration internet java javascript  
 library linux live mac media microsoft mp3 music network news online opensource osx  
 photo photography photos photoshop php portfolio productivity  
 programming python rails reference research resources rss ruby search  
 security social software tech technology tips tool tools toread tutorial  
 tutorials ubuntu video web web2.0 webdesign webdev windows xml

**Fig. 2.** Example of a tag-cloud. A tag-cloud is a common way to visualize tags belonging to a collaborative tagging system. Here, the font size of each tag is proportional to the logarithm of its frequency of appearance within the folksonomy.

(see Table 1). Fig. 3 graphically illustrates the associations between tags and posts, and Fig. 4 reports the frequency-rank distributions for the tags cooccurring with a few selected ones. The high-rank tail of the experimental curves displays a power-law behavior, signature of an emergent hierarchical structure, corresponding to a generalized Zipf's law (9) with an exponent between 1 and 2. Because power laws are the standard signature of self-organization and of human activity (10–12), the presence of a power-law tail is not surprising. The observed value of the exponent, however, deserves further investigation because the mechanisms usually invoked to explain Zipf's law and its generalizations (13) do not look very realistic for the case at hand, and a mechanism grounded on experimental data should be sought.

Moreover, the low-rank part of the frequency-rank curves exhibits a flattening typically not observed in systems strictly obeying Zipf's law. Several aspects of the underlying complex dynamics may be responsible for this feature: on the one hand, this behavior points to the existence of semantically equivalent and possibly competing high-frequency tags (e.g., *blog* and *blogs*). More importantly, this flattening behavior may be ascribed to an underlying hierarchical organization of tags cooccurring with the one we single out: more general tags (semantically speaking) will tend to cooccur with a larger number of other tags. In this scenario, we expect a shallower behavior for tags cooccurring with generic tags (e.g., *blog*) and a steeper behavior for semantically narrow tags (e.g., *ajax*, see also Fig. 3). To better probe the validity of this interpretation, we investigate the cooccurrence relationship that links high-rank tags, lying well

**Table 1. Statistics of the data sets used for the cooccurrence analysis**

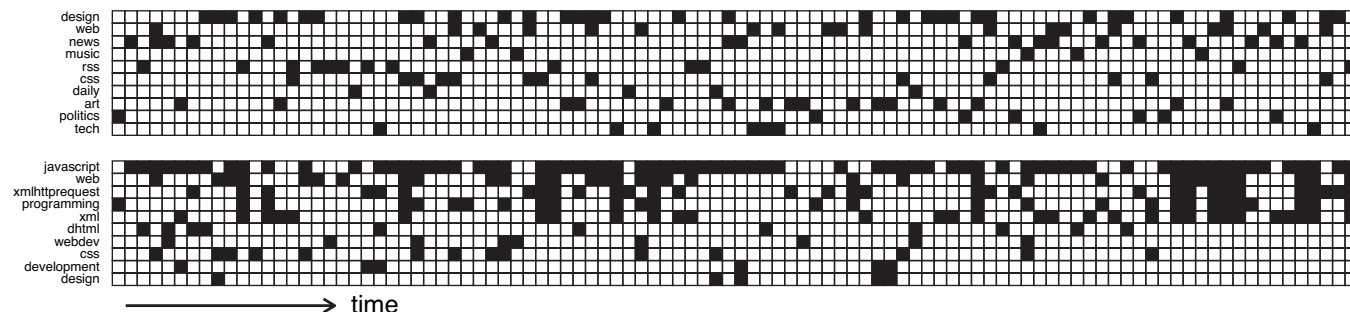
Tag	No. of posts	No. of tags	No. of distinct tags	No. of resources
Blog	37,974	24,171	10,617	16,990
Ajax	33,140	108,181	4,141	2,995
Xml	24,249	108,013	6,035	7,364
H5N1	981	5,185	241	969

For each tag, we report the number of posts marked with that tag, the number of total and distinct tags cooccurring with it, and the corresponding number of resources. The data were retrieved during May of 2005.

within the power-law tail, with low-rank tags located in the shallow part of the distribution. Our observations [see [supporting information \(SI\)](#)] point in the direction of a nontrivial hierarchical organization emerging out of the collective tagging activity, with each low-rank tag leading its own hierarchy of semantically related higher-rank tags, and all such hierarchies merging into the overall power-law tail.

**A Yule–Simon Model with Long-Term Memory.** We now aim at gaining a deeper insight into the phenomenology reported above. To model the observed frequency-rank behavior for the full range of ranking values, we introduce a new version of the “rich-get-richer” Yule–Simon's stochastic model (14, 15) by enhancing it with a fat-tailed memory kernel. The original model can be described as the construction of a text from scratch. At each discrete time step one word is appended to the text: with probability  $p$  the appended word is a new word, never occurred before, whereas with probability  $1 - p$  one word is copied from the existing text, choosing it with a probability proportional to its current frequency of occurrence. This simple process yields frequency-rank distributions that display a power-law tail with exponent  $\alpha = 1 - p$ , lower than the exponents we observe in actual data. This happens because the Yule–Simon process has no notion of “aging,” i.e., all positions within the text are regarded as identical.

In our construction, we moved from the observation that actual users are exposed in principle to all of the tags stored in the system (like in the original Yule–Simon model) but the way in which they choose among them, when tagging a new resource, is far from being uniform in time (see also refs. 16 and 17). It seems more realistic to assume that users tend to apply recently added tags more frequently than old ones, according to a memory kernel which might be highly skewed. Indeed, recent



**Fig. 3.** Tagging activity: a time-ordered sequence of tagging events is graphically rendered by marking the tags cooccurring with *blog* (Upper) or *ajax* (Lower) in an experimental sequence of posts on *del.icio.us*. In each panel, columns represent single tagging events (posts) and rows correspond to the 10 most frequent tags cooccurring with either *blog* (Upper) or *ajax* (Lower). One hundred tagging events are shown in each panel, temporally ordered from left to right. Only posts involving at least one of the 10 top-ranked tags are shown. For each tagging event (column), a filled cell marks the presence of the tag in the corresponding row, whereas an empty cell indicates its absence. A qualitative difference between *blog* (Upper) and *ajax* (Lower) is clearly visible, where a higher density at low-rank tags characterizes the semantically narrower *ajax* term. This corresponds to the steeper low-rank behavior observed in the frequency-rank plot for *ajax* (Fig. 4).





time by  $x$  steps with a probability  $Q_t(x)$  that decays as a power law,  $Q_t(x) = a(t)/(x + \tau)$  (see Fig. 6).  $a(t)$  is a normalization factor, and  $\tau$  is a characteristic time scale over which recently added words have comparable probabilities. Fig. 4 shows the excellent agreement between the experimental data and the numerical predictions of our Yule–Simon model with long-term memory. Our model, unsurprisingly, also reproduces the temporal correlation behavior observed in real data.

The interpretation of  $\tau$  (similar to that of the  $\delta$  parameter introduced above for tag-tag correlations) is related to the number of equivalent top-ranked tags perceived by users as semantically independent (see SI). In our model, in fact, the average user is exposed to a few roughly equivalent top-ranked tags and this is translated mathematically into a low-rank cutoff of the power law, i.e., the observed low-rank flattening.

Fitting the parameters of the model to match its predictions (obtained by computer simulation) against the experimental data, we obtain an excellent agreement for all of the frequency-rank curves we measured, as shown in Fig. 4. This is a clear indication that the behavior encoded in our simple model is able to capture some key features of the tagging activity. The parameter  $\tau$  controls the number of top-ranked tags that are allowed to cooccur with comparable frequencies, so that it can be interpreted as a measure of the “semantic breadth” of a tag. This picture is consistent with the fact that the fitted value of  $\tau$  obtained for *blog* (a rather generic tag) is larger than the one needed for *ajax* (a pretty specific tag). It is worth remarking that, despite the agreement between the experimental data and our model predictions, our simple modeling is an attempt toward the modeling of user behaviors, not meant to be neither unique or exclusive of other generative models (10, 11).

## Discussion and Conclusions

Uncovering the mechanisms governing the emergence of shared categorizations or vocabularies in absence of global coordination is a key problem with significant scientific and technological potential. Collaborative tagging provides a precious opportunity to both analyze the emergence of shared conventions and inspire the design of large (human or artificial) agent systems. Here we report a statistical analysis of tagging activity in a popular social bookmarking system, and introduce a simple stochastic model of user behavior which is able to reproduce the measured cooccurrence properties to a surprisingly level of accuracy. Our results suggest that users of collaborative tagging systems share universal behaviors that, despite the intricacies of personal categorization, tagging procedures, and user interactions, appear to follow simple activity patterns. In addition to the findings reported and discussed in this paper, our approach constitutes a

starting point upon which more cognitively informed studies can be based, with the final goal of understanding and engineering the semiotic dynamics of online social systems.

## Experimental Data

Our analysis focuses on *Del.icio.us*, for several reasons: (i) it was the first system to deploy the ideas and technologies of collaborative tagging, and the paradigmatic character it acquired makes it a natural starting point for any quantitative study. (ii) Because of its popularity, it has a large community of active users and comprises a precious body of raw data on the static and dynamical properties of a folksonomy. (iii) It is a broad folksonomy (as Thomas Vander Wal argued, [www.personalinfo-cloud.com/2005/02/explaining\\_and\\_.html](http://www.personalinfo-cloud.com/2005/02/explaining_and_.html)), and single tagging events (posts) retain their identity and can be individually retrieved. This affords unimpeded access to the “microscopic” dynamics of collaborative tagging, providing the opportunity to make contact between emergent behaviors and low-level dynamics. It also allows to define and measure the multiplicity (or frequency) of tags in the context of a single resource. Contrary to this, popular sites falling in the narrow folksonomy class (*Flickr*, for example) foster a different model of user interaction, where tags are mostly applied by the content creator, no notion of tag multiplicity is available in the context of a resource, and no access is given to the raw sequence of tagging events.

On studying *Del.icio.us*, we adopt a tag-centric view of the system; i.e., we investigate the evolving relationship between a given tag and the set of tags that cooccur with it. In line with our focus on semiotic dynamics, we factor out the detailed identity of the users involved in the process, and only deal with streams of tagging events and their statistical properties. To perform automated data collection of raw data we use a custom web (HTTP) client that connects to *Del.icio.us* and navigates the system’s interface as an ordinary user would do, extracting the relevant metadata and storing it for further postprocessing. *Del.icio.us* allows the user to browse its content by tag: our client requests the web page associated with the tag under study and uses an HTML parser to extract the post information (user, resource, tags, time stamp) from the returned HTML code. Fig. 3 graphically depicts the raw data we gather, for the case of two popular tags on *Del.icio.us*. Table 1 describes the data sets we used for the present analysis.

We thank A. Baronchelli, A. Baldassarri, and V. Servedio for many interesting discussions and suggestions. This research has been partly supported by the TAGora project funded by the Future and Emerging Technologies program (IST-FET) of the European Commission under the European Union RD contract IST-034721.

- Mathes A (2004) *Computer Mediated Commun*, LIS590CMC.
- Hammond T, Hannay T, Lund B, Scott J (2005) *D-Lib Magazine* 11:www.dlib.org/dlib/april05/hammond/04hammond.html.
- Golder S, Huberman BA (2006) *J Information Sci* 32:198–208.
- Steels L, Kaplan F (1999) *Lect Notes Artificial Intell* 1674:679–688.
- Steels L (2006) *IEEE Intelligent Syst* 21:32–38.
- Ke J, Minett JW, Ching-Pong A, Wang WS-Y (2002) *Complexity* 7:41–54.
- Nowak MA, Komarova NL, Niyogy P (2002) *Nature* 417:611–617.
- Kirby S (2002) *Artificial Life* 8:182–215.
- Zipf GK (1949) *Human Behavior and the Principle of Least Effort* (Addison-Wesley, Cambridge, MA).
- Mitzenmacher M (2004) *Internet Math* 1:226–251.
- Newman MEJ (2005) *Contemporary Phys* 46:323–351.
- Barabasi A-L (2005) *Nature* 435:207–211.
- Ferrer Cancho R, Servedio VDP (2005) *Glottometrics* 11:1.
- Yule GU (1925) *Philos Trans R Soc London B* 213:21–87.
- Simon HA (1955) *Biometrika* 42:425–440.
- Zanette DH, Montemurro MA (2005) *J Quant Ling* 12:29–40.
- Dorogovtsev SN, Mendes JFF (2000) *Phys Rev E* 62:1842–1845.
- Anderson JR (2000) *Cognitive Psychology and its Implications* (Worth, New York), 5th Ed.