

# Comparing the hierarchy of keywords in on-line news portals

Gergely Tibély<sup>1,\*</sup>, David Sousa-Rodrigues<sup>2</sup>, Péter Pollner<sup>3</sup>  
and Gergely Palla<sup>3</sup>

<sup>1</sup>Dept. of Biological Physics, Eötvös University, H-1117 Budapest, Hungary

\*tibelyg@hal.elte.hu

<sup>2</sup>The Design Group, Faculty of Maths, Computing and Technology  
The Open University, Walton Hall, Milton Keynes, MK7 6AA United Kingdom

<sup>3</sup>MTA-ELTE Statistical and Biological Physics Research Group,  
Hungarian Academy of Sciences, H-1117 Budapest, Hungary,

May 28, 2020

## Abstract

The tagging of on-line content with informative keywords is a widespread phenomenon from scientific article repositories through blogs to on-line news portals. In most of the cases, the tags on a given item are free words chosen by the authors independently. Therefore, relations among keywords in a collection of news items is unknown. However, in most cases the topics and concepts described by these keywords are forming a latent hierarchy, with the more general topics and categories at the top, and more specialised ones at the bottom. Here we apply a recent, cooccurrence-based tag hierarchy extraction method to sets of keywords obtained from four different on-line news portals. The resulting hierarchies show substantial differences not just in the topics rendered as important (being at the top of the hierarchy) or of less interest (categorised low in the hierarchy), but also in the underlying network structure. This reveals discrepancies between the plausible keyword association frameworks in the studied news portals.

## Introduction

Hierarchical organisation is a widespread phenomenon in nature and society. Signs of hierarchy were recorded in various animal flocks [1, 2, 3, 4], in social interactions [5, 6, 7], in urban planning [8, 9], in ecological systems [10, 11] and in evolution [12, 13]. Since a natural representation of hierarchies is given by directed acyclic graphs, hierarchical organisation became a very relevant concept also in complex network theory [14, 15, 16, 17, 18, 19, 20, 21, 22].

The association of tags to various on-line contents have become ubiquitous, as various tags may indicate the topic of news-portal feeds and blog post, the genre of films or music records on file sharing portals, or the kind of goods offered in Web stores. These tags usually serve as keywords, providing a rough description of the given entity, helping the users in a fast decision whether the given article, film, etc. is of interest or not. Keywords, categories, classes, etc. are also used in e.g., library classification systems and biological classification for helping the search and browsing amongst a large number of objects. In the latter cases the involved entities are categorised hierarchically, with a set of narrower or broader categories building up a tree-like structure composed of “is a subcategory of” type relations. In contrast, the tags appearing in on-line platforms are usually free words chosen by the author or owner of the given object, and tags are almost never organized into a pre-defined hierarchy of categories and sub-categories [23, 24, 25, 26]. Furthermore, in many tagging systems like Flickr, CiteUlike or Delicious the tagging process is collaborative, as in principle an unlimited number of users can tag photos, Web pages, etc., with free words [27, 28, 29]. In order to highlight this collaborative nature, the arising set of free tags and associated objects are often referred to as folksonomies. Since the tagging actions involve user-tag-object triplets, a natural representations of these systems is given by hypergraphs [28, 30, 31, 32, 33], where the hyperedges connect more than two nodes together.

An interesting problem related to free tagging is to extract a hierarchy between the tags based on their co-occurrences on the tagged items [34, 35, 36, 37, 38]. The basic motivation is that the way users think about objects presumably has some built in hierarchy, e.g., “pigeon” is usually considered as a special case of “bird”. Revealing the hidden hierarchy between tags in a folksonomy or in a tagging system in general can significantly help broadening or narrowing the scope of search in the system, give recommendation about yet unvisited objects to the user, or help the categorisation of newly appearing objects [33, 39]. Here we apply an improved version of a recent tag hierarchy extraction method [38] to keywords associated to on-line articles, collected from the portals of Spiegel Online, The Guardian, The New York Times and The Australian. The obtained hierarchies show very interesting differences, indicating that the methods for choosing keywords

are based on rather different principles in the studied journals.

The structure of the paper is as follows: in , the principles of hierarchy construction and comparison are shortly described. In the empirical datasets are presented, in we compare the obtained hierarchies with each other, and in results are discussed.

## Methods

### Hierarchy construction

We employ an upgraded version of a recent method [38], which is based on two assumptions:

- tags positioned high in the hierarchy also have high centrality values in the tag-tag coappearance graph,
- parent-child pairs coappear more frequently than expected from pure chance.

According to the first assumption, the algorithm orders the tags by their centrality, then, for each tag (which become child) the parent candidates are collected. All tags with higher centrality are parent candidates of the child tag. Candidate parents are assigned a score, indicating the probability of the observed number of co-occurrences according to a random null-model. Using the second assumption, the final parent is the candidate with the highest score sum, where the sum runs over all descendants of the child tag. Note, that the algorithm builds up the hierarchy bottom up, starting from the leaves with lowest centrality. The full detailed description of the currently used version of the algorithm involving a couple of improvements is given in the Supplementary Information.

### Similarity of hierarchies

Hierarchies are frequently represented by Directed Acyclic Graphs (DAGs), in which directed cycles are forbidden. However, children are allowed to have more than one parent in general. For simplicity, we have restricted the number of parents to one in the present analysis. A natural idea for comparing two DAGs is to compare the hierarchical relations, i.e., the sets of ancestor-descendant relationships [38, 40, 41, 42, 43]. Here we adopt the approach proposed in Ref. [38], defining a similarity measure based on mutual information. We note that mutual information plays a central role also in the comparison method introduced in [44] for the related, but separate problem of comparing hierarchical community structures, (where only the lowest-level nodes in DAG actually exist in the input data-set). The DAG similarity measure we use can be formulated as follows [38]

$$I_{\alpha,\beta} = \frac{2 \sum_{x=1}^{N_{\alpha\beta}} |d_{\alpha}(x) \cap d_{\beta}(x)| \cdot \ln \left( \frac{|d_{\alpha}(x) \cap d_{\beta}(x)|(N-1)}{|d_{\alpha}(x)| \cdot |d_{\beta}(x)|} \right)}{\sum_{x=1}^{N_{\alpha}} |d_{\alpha}(x)| \ln \left( \frac{|d_{\alpha}(x)|}{N-1} \right) + \sum_{x=1}^{N_{\beta}} |d_{\beta}(x)| \ln \left( \frac{|d_{\beta}(x)|}{N-1} \right)} \quad (1)$$

where  $\alpha$  and  $\beta$  are two DAGs, having  $N_{\alpha}$  and  $N_{\beta}$  tags from which  $N_{\alpha\beta}$  are common, and  $d_{\alpha}(x)$  is the set of descendants of  $x$  in DAG  $\alpha$ . Equation 1 is 0 for independent DAGs and 1 for identical ones.

A further very closely related similarity measure that turned out to be useful in previous studies is given by the linearised mutual information (LMI) [38], based on the fraction of links that have to be rewired in a randomisation procedure on  $\alpha$  leading to a hierarchy  $\alpha_{\text{rand}}$  with the same NMI when compared to  $\alpha$  as the  $I_{\alpha,\beta}$ . The formal definition of this measure is given as follows. Let  $I(f)$  denote the average NMI obtained for a fraction of  $f$  randomly rewired links,  $I(f) = \langle I_{\text{original,rand}} \rangle_f$ . By projecting the NMI of the empirical case,  $I_e$ , to the  $f$  axis using this function as

$$f^* = I^{-1}(I_e), \quad (2)$$

we receive the fraction of randomly chosen links to be rewired in the empirical case for obtaining a randomized hierarchy with the same NMI. Based on that we define the linearised mutual information, (LMI) as

$$I_{\text{lin}} = 1 - f^* = 1 - I^{-1}(I_e) \quad (3)$$

This quality measure corresponds to the fraction of unchanged links in a random link rewiring process, resulting in a hierarchy with the same NMI as the empirical value. (The reason for calling it “linearised” is that equation 3 is actually projecting  $I_e$  to the linear  $1 - f$  curve).

## Data

We analyse four tagged datasets, obtained from online news portals. They contain tagged news items, covering a more than 2 years long time window, in the same period. The four sources are: Spiegel Online, The Guardian, The New York Times and The Australian.

### General observations

There are a few observations which hold for all four datasets. For example, very long tags exist, more like headlines (“**Muntazer al-Zaidi: the Iraqi shoe thrower**”). Some of the tags form “frozen” cliques in the coappearance network, where each member of such a clique appear only together with the other members of the clique, e.g., **Haiti** and **Haiti Earthquake Disaster 2010**, **Diana** and **Princess of Wales**. Since members of a large

clique have large centrality values, such tags will be placed to unwanted high positions by the first step of the hierarchy construction algorithm. Therefore we have considered such “frozen” cliques as single tags, which fits better to the assumed usage of tags.

Some concepts are represented by two or more tags, where the same idea is expressed with different, but synonymous words, e.g., **Art** and **Arts**. These were left as observed, unless explicitly stated otherwise. Another problem is posed by the occurrence of very rare tags, that are usually names.

In order to avoid misleading results due to the above observed problems, we have prefiltered the tags by requiring that each tag pair in the coappearance network has to occur on at least  $r$  news items. The  $r = 1$  case corresponds to skipping the prefiltering. We set  $r$  to its optimal value for each dataset by keeping the number of tags as high as possible and minimizing the number of misleading tags described above. Finally we note, that temporarily important topics can produce unexpected co-occurrences (e.g., “Japan”  $\rightarrow$  “Fukushima Nuclear Catastrophe”  $\rightarrow$  “Nuclear Power”).

### **Spiegel Online**

The dataset is from April 2011 to January 2013. It contains 4802 news items and 388 tags. For the pre-filtering, minimum 1 common news item for each tag pair (i.e., no filtering) seems to be a good tradeoff between noise reduction and info loss. The dataset looks very well organised (e.g., there are only 400 tags, general tags are used consistently, and there are only a few duplicated tags, long tags or “frozen” cliques).

### **The Guardian**

The dataset is from November 2009 to January 2013, containing 55835 news items and 6797 tags. Pre-filtering needs minimum 3 news items (removes 2530 tags and 61 news items). Here we found several ad hoc tags (mostly names), that were used only once or a handful of times. We found synonymous tags, e.g., “Middle East and North Africa” and “Middle East”, that will appear as two local roots of two branches in the DAG. These branches correspond to the same topic, thus divide the related tags between them.

### **The New York Times**

The dataset reaches from November 2010 to January 2013. It contains 35736 news items and 23009 tags. Cliques are a huge problem here. There are 2902 ones, collapsing them removes about 6000 tags. Several cliques appear on numerous objects, therefore the minimum news item-filtering does not solve the problem automatically. Cliques also reach very large sizes: there is a 809-tag clique (may contain much more characters than a news item itself); after the minimum news items filtering, the largest one still consists of 44 tags – as follows from the definition of cliques, these tags appear strictly together on each object. For the pre-filtering, minimum 5

news items were required, leaving finally 2981 tags (out of 23009). News items were much less affected, 31184 out of 35736 remained.

### **The Australian**

Data is from December 2009 to January 2013. It contains 31501 news items and 79054 tags – thus, there are much more tags than news items. Cliques are present, but have only 1-2 objects, so it is not a serious problem, the pre-filtering can solve it. Multiple synonyms occur on the same object very often – e.g., "Economist\_Paul\_Samuelson" "Paul\_A.\_Samuelson" "Paul\_Samuelson". Another example is the set of synonyms for Barack Obama, which are: Barack\_Obama, BARACK\_Obama, Obama, PRESIDENT\_Barack\_Obama, President\_Barack\_Obama, President\_Obama, US\_PRESIDENT\_Barack\_Obama, US\_President\_Barack, US\_President\_Barack\_Obama, barack\_obama. Pre-filtering with minimum 5 news items leaves 1673 tags out of 79504. The news items are reduced from 31501 to 10550. The tags have relatively few objects, and not only due to the large number of very infrequent tags, e.g., even the prime minister has only 900 objects. Although there are very general tags like `community`, `committee` or `claim`, most of tags are very specific, almost tailored for one object, e.g., `rebels_storm_Gaddafi_compound`.

## **Results**

We analysed the tag hierarchies obtained from an improved version of "algorithm B" published in Ref. [38]; a brief description of the idea of the method can be found in , the full details of the used algorithm are given in the Supplementary Information. In first we summarise the most important properties of the individual hierarchies corresponding to the different news portals, which is followed by the pairwise comparisons in . Finally, in we examine the overall quality of the hierarchies from different aspects.

### **Analysis of the individual tag hierarchies**

**Spiegel Online** The constructed DAG consists of 1 connected component. Most of the tags are under 3 branches: "World", "Europe", "Germany". A visualisation of the DAG is shown on Fig. 1. The Spiegel DAG seems to be somewhat concerned with immigrants and integration, they have a branch containing 3.9% of the tags, similarly to Australian's 4.4%, and in contrast to 0.1% and 0.7% of Guardian and NYT (note that the latest data come from January 2013, well before the beginning of the recent migrant crisis).

**The Guardian** The overall structure of the DAG is quite well organised, the top 2-3 levels are very impressive. The DAG consists of four similarly-sized connected components: "UK news", "World news", "Culture", "Sport", although the tags "World news" and "UK news" are in isolated components, they are not completely mutually exclusive, e.g., both of them appear on the news items of "Defence policy". Note that while the components' top tags correspond well to the menu items on the journal's website, they are

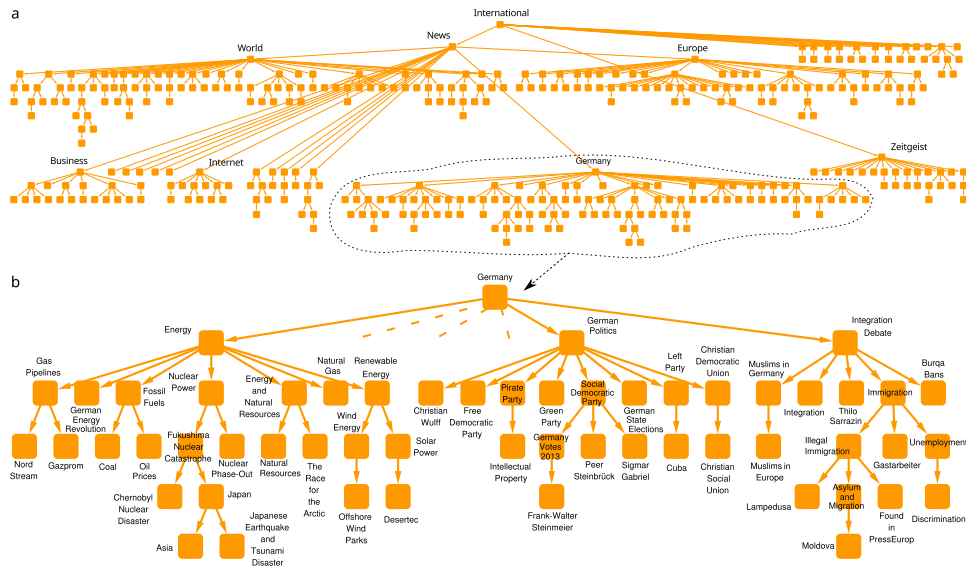


Figure 1: Overview of the Spiegel DAG (top), and one part enlarged (bottom). The DAG is broken into two lines in the top figure to fit the whole graph in the available width. On the bottom figure, dashed lines indicate descendants which are not shown.

placed totally automatically by the DAG construction algorithm. Visualisation is omitted due to the relatively large size of the DAG, however, a smaller sample is shown in the Supplementary Information.

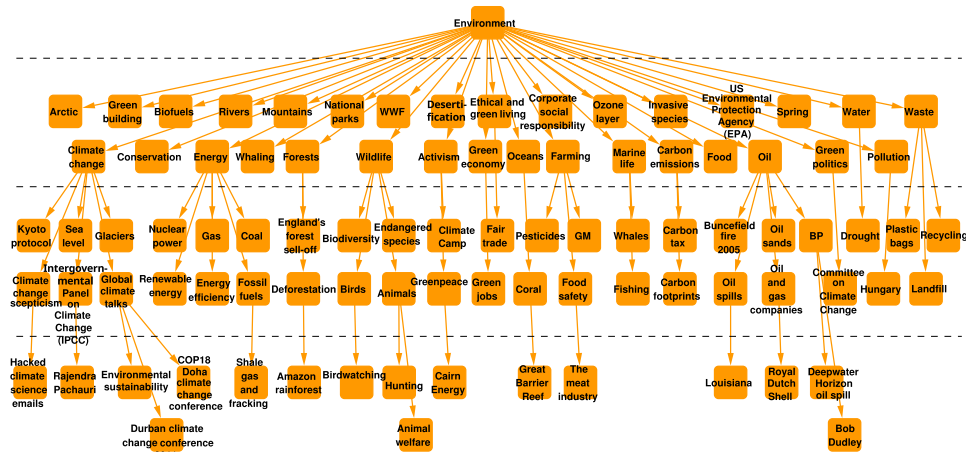


Figure 2: Part of the Guardian's Environment branch, in the component World news. Hierarchical levels are separated by dashed lines.

**The New York Times** Here we found numerous duplicated branches in the constructed DAG (e.g., for research, television, education, medicine, defence and military forces). This indicates that for these topics, two distinct sets of tags were used in parallel. The DAG is much less organised than that of the Spiegel and of the Guardian. There are 31 isolated components, most of them correspond to one theme (e.g. "Baseball"). The sizes of the components varies from 898 to 2, and there is a continuous range of them from the 2nd largest one (274 tags) down. There are no very general categories. Although a number of large related components exists (under the tags "Basketball", "Baseball", "Football"), these components are not collected under a general "Sport" tag. It seems as if there were no demand for using general tags. Note that there is a tag called "sports", however, it appears only on 5 news items, and it is negligible. A technical consequence is that the DAG construction algorithm does not always select the most general tags as roots, because they lack the important connections to other components. Instead, one of the more specific tags can be selected for a central position, for example, "Middle East and North Africa Unrest (2010-)" for foreign affairs, or "European Sovereign Debt Crisis (2010-)" for Europe-related tags. In other words, the centrality no longer correlates only with the generality for the top tags. Some lower-level branches end up at unexpected places, e.g., **Environment** under **Iran**. Superfluous levels appears, for example, **International Relations** under **United States International Relations**.

**The Australian** The DAG looks disorganised overall. There are about 1900 components for the 79504 tags without the pre-filtering, and about 300 components for the min. 5 news items-filtered 1673. There are no macroscopic components, the largest one's size is just 3480 (out of 79504 tags) and 165 (out of 1673 tags), which is less than 10% of the total nodes. Even the existing components look more like just bunches of more or less associated tags than small hierarchical structures.

In general, the top of the constructed DAGs are much better than the bottom. This is no surprise - there is much more information for the construction algorithm at the top of the DAG.

## Pairwise comparisons

We carried out a pairwise comparison between the journals from the point of view of their content organisation. Since the audience and the interests of the journals are different, the list of tags appearing on the articles was unique for each news portal. Therefore, before actually comparing the tag hierarchies, first we needed to create a common tag set for each pair of journals. In a number of cases, finding the corresponding tag pairs went beyond a simple string matching and was based on semantic matching, e.g., "Fossil fuels" (Guardian) was matched with "Oil



(Petroleum) and Gasoline” (NYT). The size of the reduced common tag sets were 252 (Spiegel-Guardian), 217 (Spiegel-NYT), 985 (Guardian-NYT), 93 (Australian-Spiegel), 278 (Australian-Guardian), 274 (Australian-NYT).

The reduced hierarchies were obtained by keeping only the common tags in the original DAGs and erasing the rest of the tags. In most cases this resulted in deletion of leafs, sub-branches, or lower parts of sub-branches from the original hierarchies. However, a small number of times this procedure erased a tag higher in a given branch while keeping other tags lower in the same branch, therefore distorting the original DAG structure in a radical way. To ensure as much similarity to the original hierarchies as possible, under these circumstances the ancestors standing higher in the branch were also kept, despite that they were not part of the common tag set, (see the SI for more details). The reduced DAGs can be found in the SI.

For each pair of journals we have computed the linearised information similarity measure described in between the reduced DAGs, the obtained values are shown in Fig. 3. According to the results Spiegel and Guardian provide the largest similarity measure, which is also supported by a number of identical or almost identical sub-branches between the two DAGs, as shown in Fig. 4. Here the background colouring of the sub-branches indicate the similarity to the corresponding (most similar) sub-branch in the other DAG.

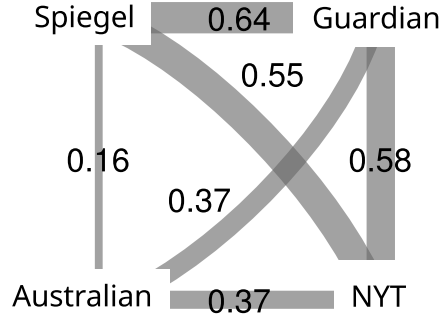


Figure 3: Similarities between the news portals’ DAGs, according to the mutual information-based linearised information similarity measure described in .

The Spiegel, the Guardian and the New York Times have an overall similar structure, as Fig. 3 shows, opposed to the Australian, which is dissimilar to all of them. Still, there are differences between the first three journals. The Guardian, compared to the Spiegel, has a level of intermediate-sized branches, e.g., **law** or **society** in **UK news**. This level is missing from the DAG of Spiegel. Their global DAG structures are shown in Fig. 4. Meanwhile, the New York Times has interestingly no **World** tag, and foreign

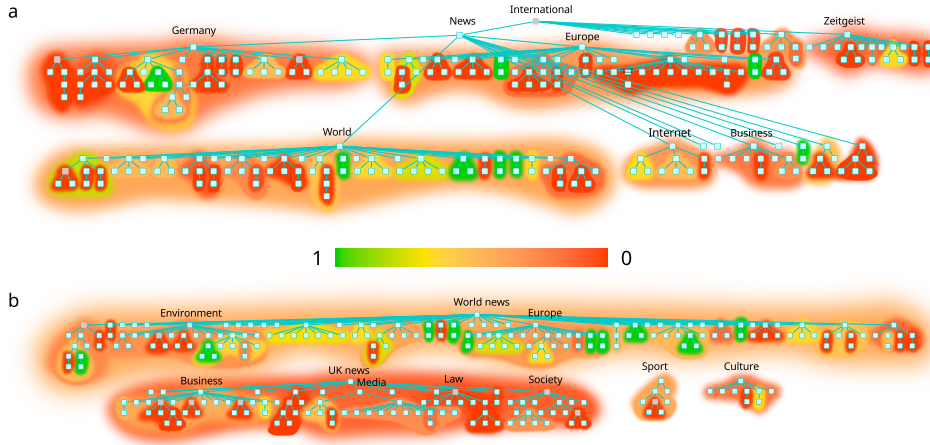


Figure 4: The Spiegel (top) and the Guardian’s (bottom) reduced DAG structures, providing the largest overall similarity in our analysis. For clarity, Spiegel’s DAG is broken into two lines. Background colours show the result obtained by applying the similarity measure given in equation 1 to the given branch and the most similar branch from the other hierarchy. Note that sub-branches on all hierarchical levels have their own colour.

countries are separated into 4 different branches, in 3 components (see the SI for more details). Although the linearised information similarity between the Guardian and the New York Times is somewhat lower, they also have a few quite similar branches; a prominent example is shown in Fig. 5.

### Statistical properties of the overall hierarchy structures

According to the results presented in the previous sections the tag hierarchies obtained for the studied journals show strong differences. Here we examine to what extent does their overall structure follow a few simple intuitive requirements that can be formulated for a well organised tag hierarchy.

**Correlations with Google News.** One of the basic properties of a well organised hierarchy is that frequent, more general tags are expected to be higher compared to rare, specific tags. In order to examine the obtained hierarchies from this perspective we compared the centrality score of the tags in the tag co-occurrence network (determining their position in the hierarchy) with their number of hits provided by Google News. For each pair of tags with a significant number of co-occurrence we checked whether the difference between their centrality score and the difference between their number of hits in Google News have the same or the opposite sign. If the signs of the differences match for the majority of the tag pairs, then we can

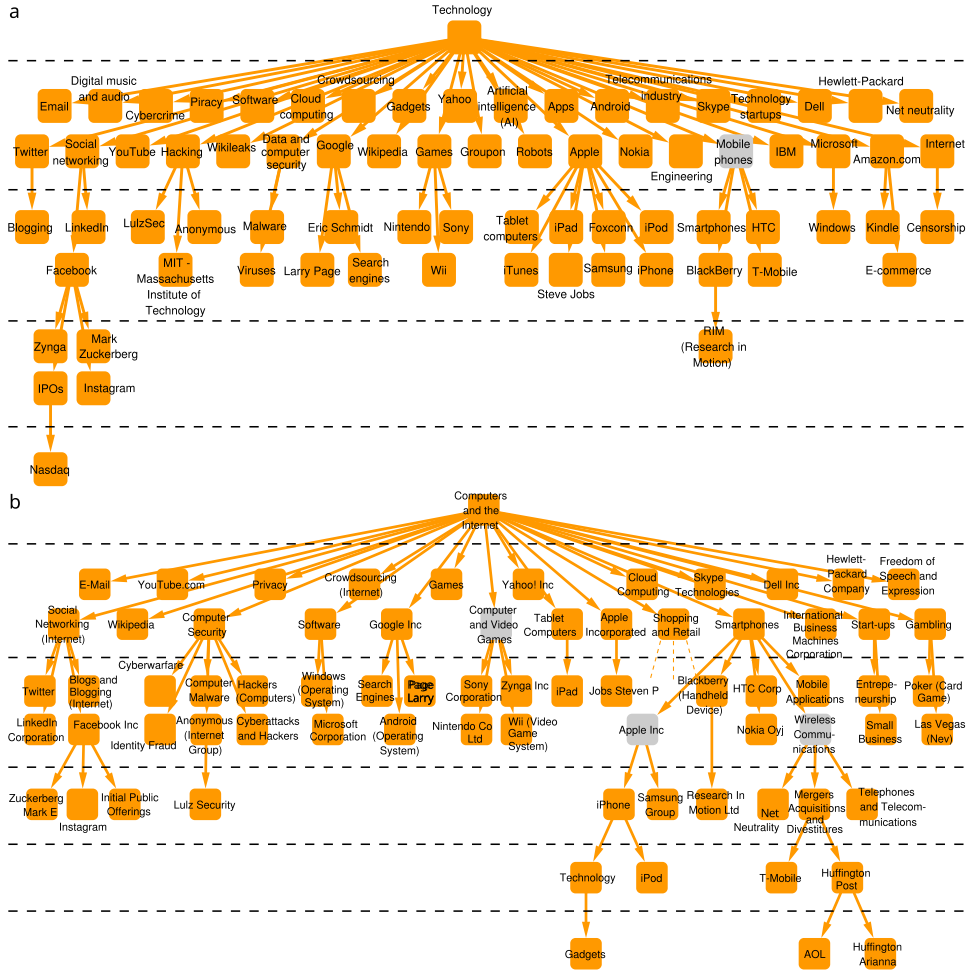


Figure 5: Guardian’s **Technology** branch (top) and New York Times’ **Computers and the Internet** component (bottom). Hierarchical levels are separated by dashed lines. Grey tags do not appear in both DAGs, however, they connect branches containing common tags.

assume that the structure of the hierarchy is consistent with word frequencies of English news texts around the world.

In Table 1 we show the relative frequency of the cases, where the differences have the opposite sign, calculated for tag pairs co-appearing in statistically significant numbers. If tags are assigned to articles absolutely at random, the result would correspond to a 0.5 inversion rate, i.e., half of the coappearing tag pairs would have similar centrality and frequency ordering. According to Table 1, the Spiegel and the Guardian data sets provide the best correspondence between tag frequency and centrality, with only a few percent difference in their score. They are followed by the New York Times,

Table 1: Ratios of inversions between centralities and real-world occurrence frequencies, calculated for tag pairs coappearing in statistically significant numbers. Totally random case corresponds to 0.5.

dataset	ratio of inversions
Spiegel	0.19
Guardian	0.21
New York Times	0.31
Australian	0.44

and finally, the Australian has a score close to the random case. Although the Google News data may be somewhat different from a fictitious collection word usage of all English speaking journalist, the results in Table 1 show a quite clear-cut picture, which also corresponds well to the results of other comparisons.

**Geometrical properties of the hierarchies.** In this section we focus on the geometrical properties of the tag hierarchies from the perspective of whether their structure is helping navigability and search. First we examine the fragmentation of the DAGs, which we can quantify by first introducing the average size of the component of a randomly chosen tag given by,

$$\tilde{s} = \frac{\sum_i^{\text{tags}} s_i}{N} \quad (4)$$

where  $s_i$  is the size of the component containing tag  $i$  and  $N$  is the total number of tags. Based on  $\tilde{s}$  we can calculate the expected lowest hierarchy level  $l$  on which the top node of a branch of size  $\tilde{s}$  would appear in a balanced  $k$ -ary tree of size  $N$ . In such a tree any branch can contain at most half of the tags of its mother branch, thus we define  $l$  as

$$\begin{aligned} l &= \lceil \log_2 N / \tilde{s} \rceil, & \tilde{s} < N \\ l &= 1, & \tilde{s} = N \end{aligned} \quad (5)$$

where  $\lceil x \rceil$  denotes the ceiling function of  $x$ . The value of  $l$  becomes high for strongly fragmented tag hierarchies consisting of many small isolated components, where the navigability of the hierarchy is low. The results for  $\tilde{s}$  and  $l$  are summarised in Table 2. The tag hierarchy obtained for Spiegel (consisting of a single component) provides the lowest  $l$  value, followed by Guardian and New York times. Apparently, the DAG of Australian is showing a very fragmented structure with  $l = 6$ .

Another important question is whether branch sizes are balanced or not in the hierarchies. A well-balanced hierarchy is expected to have at

Table 2: Characteristic level showing the highest level of an idealised hierarchy to which an average connected component corresponds.

dataset	$\tilde{s}$	$N$	$l$
Spiegel	388	388	1
Guardian	1338.7	4263	2
New York Times	384.2	2945	3
Australian	46.2	1487	6

least 2 but not more than  $\mathcal{O}(1)$  comparably sized branches at every nonleaf tag. We define a balancedness measure with a pair of real numbers from  $[0, 1) \times [0, 1)$  corresponding to the ratio of “giant branches” and the ratio of “dwarf branches” in order to quantify how a DAG fits to the above criterion. First, we calculate the cumulated size of the branches having a child branch which contains more than 50% of the parent branch’s tags. Second, we calculate the cumulated size of the child branches which are smaller than 10% of their parent branches. The higher threshold is motivated by the fact that a child branch above 50% is larger than all the other child branches combined. The motivation for the lower threshold is that below 10%, for equal-sized child branches, the number of child branches exceeds  $\mathcal{O}(1)$ . Other numerical threshold values might also be applied, however, for demonstrating significant phenomena the precise value of the thresholds should not be important. We normalise the sums by their maximal possible value, thus, our balancedness measure is given by

$$(R_g, R_d) = \left( \frac{\sum_g S_g}{\sum_b S_b}, \frac{\sum_d S_d}{\sum_b S_b} \right) \quad (6)$$

where  $b$  goes over all branches containing at least 2 tags,  $S_b$  is the size of branch  $b$ ,  $g$  goes over branches containing a sub-branch having more than 50% of  $g$ ’s tags, and  $S_g$  is the corresponding branch size, and  $d$  goes over sub-branches which are smaller than 10% of their parent branches with  $S_d$  being the corresponding branch size. A perfectly balanced hierarchy would have a (0,0) score and the two extremely unbalanced cases would have (1,0) for a chain and (0, 1) for a star graph. The results for  $(R_g, R_d)$  are given in Table 3.

Spiegel’s  $R_g$  is dominated by a single contribution. The global root, **International** has a branch containing almost the whole DAG under **News**. Most of  $R_d$  comes from small branches, although there are a few exceptions. In the Guardian DAG, dwarf sub-branches are common, due to the huge size of the components which dwarf several branches, as well as to nearly star-shaped branches, sometimes containing hundreds of leaf-tags (e.g., **Film**, **Music**). For the NYT, contrary to the Guardian,  $R_g$  is much larger than  $R_d$ .

Table 3: Ratios of giant and dwarf branches among all branches, size-weighted.

dataset	$R_g$	$R_d$
Spiegel	0.32	0.22
Guardian	0.10	0.42
New York Times	0.42	0.22
Australian	0.26	0.17

Two important reasons are misplacing a number of branches and letting less general tags getting high centralities. Since the Australian DAG has quite limited structure inside the numerous small components,  $R_g$  and  $R_d$  are not very informative measures here. However, the tiny components seem to be well balanced.

Further analysis of the DAGs can be found in the SI.

## Discussion

We studied the hierarchy of keywords associated to news articles in four different on-line news portals. The datasets contain various artefacts, such as long and complex keywords, “frozen” cliques of exclusively coappearing tags, synonyms or very rare and specific tags. Nonetheless, it was possible for the construction method to obtain very reasonable DAGs from the data. The identification of frozen cliques might also be applied by disambiguation techniques, to identify cliques of equivalent semantic meaning, used in the field of Natural Language Processing. The constructed DAGs suggest that the tags appearing in the different news portals are organised to different degrees. Our analysis revealed that Guardian has an extra intermediate level of organisation at certain locations. A further very interesting result is that the number of connected components in the DAGs conveys information about the extent of organisation in the data: the Spiegel and Guardian have  $\mathcal{O}(1)$  components and are quite organised, the New York Times has a few dozen components and breaks the world into independent pieces, and the Australian has  $\mathcal{O}(100)$  components which are barely informative at all.

A similar picture was emerging from the comparison between the frequencies of tags in Google News and their centrality score in the tag-tag co-appearance graphs. The correlation was quite strong in case of the Spiegel and the Guardian, medium for the New York Times, and almost equivalent to the totally random case for the Australian. A more detailed characterisation of the DAGs can be obtained by quantifying the extents of too large and too small sub-branches. Although being a geometry-based analysis, it can also identify problems with tag functions, like a non-comprehensive set of intermediate-level branches in the Guardian, or misplaced branches in the

New York Times.

In summary, the following picture is arising from the different analyses we carried out: the Spiegel and Guardian datasets are quite well-organised, the New York Times is significantly less but still has relevant hierarchical structure, and the Australian is close to being random, from a hierarchical point of view. The consistency of the results is encouraging, and suggests that the measures used are useful in the quantification and comparison of datasets from the aspect of hierarchical organisation.

## References

- [1] Goessmann, C., Hemelrijk, C. & Huber, R. The formation and maintenance of crayfish hierarchies: behavioral and self-structuring properties. *Behav. Ecol. Sociobiol.* **48**, 418–428 (2000).
- [2] Nagy, M., Ákos, Z., Biro, D. & Vicsek, T. Hierarchical group dynamics in pigeon flocks. *Nature* **464**, 890–893 (2010).
- [3] Nagy, M. et al. Context-dependent hierarchies in pigeons. *Proc. Natl. Acad. Sci. USA* **110**, 13049–13054 (2013).
- [4] Fushing, H., McAssey, M. P., Beisner, B. & McCowan, B. Ranking network of captive rhesus macaque society: a sophisticated corporative kingdom. *PLoS ONE* **6**, e17817 (2011).
- [5] Guimerà, R., Danon, L., Díaz-Guilera, A., Giralt, F. & Arenas, A. Self-similar community structure in a network of human interactions. *Phys. Rev. E* **68**, 065103 (2003).
- [6] Pollner, P., Palla, G. & Vicsek, T. Preferential attachment of communities: the same principle, but a higher level. *Europhys. Lett.* **73**, 478–484 (2006).
- [7] Valverde, S. & Solé, R. V. Self-organization versus hierarchy in open-source social networks. *Phys. Rev. E* **76**, 046118 (2007).
- [8] Krugman, P. R. Confronting the mystery of urban hierarchy. *J. Jpn. Int. Econ.* **10**, 399–418 (1996).
- [9] Batty, M. & Longley, P. *Fractal Cities: A Geometry Of Form And Function* (Academic, 1994).
- [10] Hirata, H. & Ulanowicz, R. Information theoretical analysis of the aggregation and hierarchical structure of ecological networks. *J. Theor. Biol.* **116**, 321–341 (1985).
- [11] Wickens, J. & Ulanowicz, R. On quantifying hierarchical connections in ecology. *J. Soc. Biol. Struct.* **11**, 369–378 (1988).

- [12] Eldredge, N. *Unfinished Synthesis: Biological Hierarchies And Modern Evolutionary Thought* (Oxford University Press, 1985).
- [13] McShea, D. W. The hierarchical structure of organisms. *Paleobiology* **27**, 405–423 (2001).
- [14] Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabási, A.-L. Hierarchical organization of modularity in metabolic networks. *Science* **297**, 1551–1555 (2002).
- [15] Trusina, A., Maslov, S., Minnhagen, P. & Sneppen, K. Hierarchy measures in complex networks. *Phys. Rev. Lett.* **92**, 178702 (2004).
- [16] Clauset, A., Moore, C. & Newman, M. E. J. Hierarchical structure and the prediction of missing links in networks. *Nature* **453**, 98–101 (2008).
- [17] Pumain, D. *Hierarchy In Natural And Social Sciences*, Vol. 3 *Methodos Series* (Springer, 2006).
- [18] Corominas-Murtra, B., Rodríguez-Caso, C., Goñi, J. & Solé, R. V. Measuring the hierarchy of feedforward networks. *Chaos* **21**, 016108 (2011).
- [19] Mones, E., Vicsek, L. & Vicsek, T. Hierarchy measure for complex networks. *PLoS ONE* **7**, e33799 (2012).
- [20] Corominas-Murtra, B., Goñi, J., Solé, R. V. & Rodríguez-Caso, C. On the origins of hierarchy in complex networks. *Proc. Natl. Acad. Sci. USA* **110**, 13316–13321 (2013).
- [21] Kaiser, M., Hilgetag, C. C. & Kötter, R. Hierarchy and dynamics of neural networks. *Front. Neuroinform.* **4**, 112 (2010).
- [22] Ma, H. W., Buer, J. & Zeng, A. P. Hierarchical structure and modules in the escherichia coli transcriptional regulatory network revealed by a new top-down approach. *BMC Bioinform.* **5**, 199 (2004).
- [23] Mika, P. Ontologies are us: a unified model of social networks and semantics in *The Semantic Web – ISWC 2005*, Vol. 3729 (eds Gil, Y., Motta, E., Benjamins, V. R. & Musen, M. A.) in *Lect. Notes Comput. Sci.*, 522–536 (Springer, 2005).
- [24] Spyns, P., Moor, A. D., Vandenbussche, J. & Meersman, R. From folk-sociologies to ontologies: how the twain meet in *On The Move To Meaningful Internet Systems 2006: CoopIS, DOA, GADA, and ODBASE*, Vol. 4275 (eds Meersman, R. & Tari, Z.) in *Lect. Notes Comput. Sci.*, 738–755 (Springer, 2006).



- [25] Voss, J. Tagging, folksonomy & co - renaissance of manual indexing? *arXiv preprint arXiv:cs/0701072v2* (2007).
- [26] Tibély, G., Pollner, P., Vicsek, T. & Palla, G. Ontologies and tag-statistics. *New J. Phys.* **14**, 053009 (2012).
- [27] Cattuto, C., Loreto, V. & Pietronero, L. Semiotic dynamics and collaborative tagging. *Proc. Natl. Acad. Sci. USA* **104**, 1461–1464 (2007).
- [28] Lambiotte R. & Ausloos, M. Collaborative tagging as a tripartite network in *Computational Science – ICCS 2006*, Vol. 3993 (eds Alexandrov, V. N., van Albada, G. D., Soot, P. M. A. & Dongarra, J.) in *Lect. Notes Comput. Sci.*, 1114–1117 (Springer, 2006).
- [29] Cattuto, C., Barrat, A., Baldassarri, A., Schehr, G. & Loreto, V. Collective dynamics of social annotation. *Proc. Natl. Acad. Sci. USA* **106**, 10511–10515 (2009).
- [30] Ghoshal, G., Zlatić, V., Caldarelli, G. & Newman, M. E. J. Random hypergraphs and their applications. *Phys. Rev. E* **79**, 066118 (2009).
- [31] Zlatić, V., Ghoshal, G. & Caldarelli, G. Hypergraph topological quantities for tagged social networks. *Phys. Rev. E* **80**, 036118 (2009).
- [32] Floeck, F., Putzke, J., Steinfels, S., Fischbach, K. & Schoder, D. Imitation and quality of tags in social bookmarking systems - collective intelligence leading to folksonomies in *Advances In Intelligent And Soft Computing*, Vol. 76 (eds Bastiaens, T. J., Baumöl, U. & Krämer, B. J.) Ch. On Collective Intelligence, 75–91 (Springer, 2011).
- [33] Lu, L. et al. Recommender systems. *Phys. Rep.* **519**, 1–49 (2012).
- [34] Heymann, P. and Garcia-Molina, H., *Collaborative creation of communal hierarchical taxonomies in social tagging systems. Stanford InfoLab Technical Report.* (2006) Available at: <http://ilpubs.stanford.edu:8090/775/?auth=basic>. (Accessed: 18th January 2012)
- [35] Plangprasopchok, A., Lerman, K. & Getoor, L. A probabilistic approach for learning folksonomies from structured data. In *Fourth ACM International Conference On Web Search And Data Mining (WSDM)*, 555–564 (2011).
- [36] Schmitz, P. Inducing ontology from flickr tags. In *Proceedings Of Collaborative Web Tagging Workshop At The 15th International Conference On World Wide Web (WWW)* (2006).

- [37] Van Damme, C., Hepp, M. & Siorpaes, K. Folksontology: an integrated approach for turning folksonomies into ontologies. *Soc. Networks* **2**, 57–70 (2007).
- [38] Tibély, G., Pollner, P., Vicsek, T. & Palla, G. Extracting tag hierarchies. *PLoS ONE* **8**, e84133 (2013).
- [39] Juszczyszyn, K., Kazienko, P. & Katarzyna, M. Personalized ontology-based recommender systems for multimedia objects in *Studies In Computational Intelligence*, Vol. 289 (eds Håkansson, A., Hartung, R. & Nguyen, N.) Ch. Agent and Multi-agent Technology for Internet and Enterprise Systems, 275–292 (Springer, 2010).
- [40] Fattore, M., Grassi, R. & Arcagni, A. Measuring structural dissimilarity between finite partial orders in *Multi-indicator Systems And Modelling In Partial Order* (eds Brüggemann, R., Carlsen, L. & Wittmann, J.), 69–84 (Springer, 2013).
- [41] Brandenburg, F. J., Gleißner, A. & Hofmeier, A. The nearest neighbor spearman footrule distance for bucket, interval, and partial orders. *J. Comb. Optim.* **26**, 310–332 (2013).
- [42] Palla, G., Tibély, G., Mones, E., Pollner, P. & Vicsek, T. Hierarchical networks of scientific journals. *Palgrave Communications* **1**, 15016 (2015).
- [43] Tibély, G., Pollner, P. & Palla, G. Partial order similarity based on mutual information. *arXiv preprint arXiv:1601.05922* (2016).
- [44] Tessone, C. J., Perotti, J. I. & Caldarelli, G. Hierarchical mutual information for the comparison of hierarchical community structures in complex networks. *Phys. Rev. E* **92**, 062825 (2015).

## Acknowledgements

Financial support of the Hungarian National Science Fund (OTKA K105447) is acknowledged.

## Author Contributions

GT, PP and GP conceived and designed the experiments. DSR provided the datasets. GT analysed the data. GT and GP wrote the manuscript. All authors reviewed the manuscript.

## Additional Information

**Competing interests:** The authors declare no competing financial interests.