

# Modelling of temporal fluctuation scaling in online news network with independent cascade model

Jan Chołoniowski<sup>a</sup>, Julian Sienkiewicz<sup>a</sup>, Gregor Leban<sup>b</sup>, Janusz A. Hołyst<sup>a,c,\*</sup>

<sup>a</sup>*Center of Excellence for Complex Systems Research, Faculty of Physics, Warsaw University of Technology, Koszykowa 75, 00-662, Warsaw, Poland*

<sup>b</sup>*Artificial Intelligence Laboratory, Jožef Stefan Institute, Jamova 39, 1000, Ljubljana, Slovenia*

<sup>c</sup>*ITMO University, 49 Kronverkskiy av., 197101, Saint Petersburg, Russia*

---

## Abstract

We show that activity of online news outlets follows a temporal fluctuation scaling law and we recover this feature using an independent cascade model augmented with a varying *hype* parameter representing a viral potential of an original article. We use the Event Registry platform to track activity of over 10,000 news outlets in 11 different topics in the course of the year 2016. Analyzing over 22,000,000 articles, we found that fluctuation scaling exponents  $\alpha$  depend on time window size  $\Delta$  in a characteristic way for all the considered topics – news outlets activities are partially synchronized for  $\Delta > 15\text{min}$  with a cross-over for  $\Delta = 1\text{day}$ . The proposed model was run on several synthetic network models as well as on a network extracted from the real data. Our approach discards timestamps as not fully reliable observables and focuses on co-occurrences of publishers in cascades of similarly phrased news items. We make use of the Event Registry news clustering feature to find correlations between content published by news outlets in order to uncover common information propagation paths in published articles and to estimate weights of edges in the independent cascade model. While the independent cascade model follows the fluctuation scaling law with a trivial exponent  $\alpha = 0.5$ , we argue that besides the topology of the underlying cooperation network a temporal clustering of articles with similar hypes is necessary to qualitatively reproduce the fluctuation scaling observed in the data.

**Keywords:** Fluctuation scaling, Complex systems, Complex networks, Online media, Agent-based modelling

Declarations of interest: none

---

\*Corresponding author.

*Email addresses:* choloniowski@if.pw.edu.pl (Jan Chołoniowski), julas@if.pw.edu.pl (Julian Sienkiewicz), gregor.leban@ijs.si (Gregor Leban), jholyst@if.pw.edu.pl (Janusz A. Hołyst)

## 1. Introduction

Rapid digitalization of our everyday life created possibilities to analyze previously inaccessible areas. Specifically, the emergence of Web 2.0 and social media (e.g., Twitter, blogosphere, Facebook) gave foundation to computational social sciences [1]. A physicists' involvement can be seen as a successful one [2, 3, 4, 5, 6, 7]; examples of the latest studies leveraging the recent advancements are quantifying and modelling emotions [8, 9, 10] and opinions [11] in online communities, network of scientific collaborations [12, 13, 14, 15], dynamics of languages [16], or information diffusion and processing [17, 18, 19, 20, 21].

The social media revolution caused major changes in media industry. The Internet turned out to be an extremely effective medium for news stories changing the modern journalism – publishing and spreading are more dynamic now, news is produced and updated continuously; also the gatekeeping function of journalists has been reduced since every person or entity can directly post messages on social media. For data scientists, the shift towards digitalized environment created an opportunity to quantitatively analyze activity and content of news outlets. One can perceive online news ecosystem as a complex network consisting of news outlets producing information themselves or mimicking/processing information observed in other sources (such as news agencies, social media, or other news outlets). Such dynamical systems often follow some kind of a fluctuation scaling law [22] which can be identified to provide conclusions about degree of units' temporal [23] or spatial [24] synchronization in different scales.

The study presented here is an analysis of 22 million articles with one of 11 keywords published in 2016 gathered by *EventRegistry.org* [25], a global media monitor. The main aim of our study is to report the temporal fluctuation scaling found in the dataset and check whether the statistical property of media activity can be described with an epidemic-like process. The observed type of the fluctuation scaling can give hints on underlying system dynamics. We show that a possible explanation of the observed fluctuation scaling could be an underlying independent cascade model taking place on a complex network. We track propagation of news stories to approximate a topology of the network. While in many situations information is propagated as a straightforward connection, like retweets, in other cases it may mutate, change its form or sentiment, or even become a mix of information from different sources. In journalism, explicitly mentioning a cited source is considered a good practice but, due to the competitiveness of the industry, it is not always met. Without reliable information about which news item was published first (as it might be a matter of seconds), whether given piece of content was produced or copied, or even whether its original source is being observed, we decided to apply natural language processing methods to meaningfully group published news items across various news outlets [26, 27]. Aggregating results across another dataset (over 14,000 events observed between 1.5.2017 and 8.5.2017 limiting to articles written in English) uncovered a content correlation network which was used to simulate the process of news spreading. We show that the independent cascade model in its common form does not provide full explanation of the fluctuation scaling observed in

the data, and postulate a specific news item feature – *hype* – which represents its intrinsic viral potential and causes the model to indicate realistic scaling exponents.

There is an abundance of recent methods to uncover network basing on the independent cascade model (typically continuous time [28, 29, 30]) which assume exactly measured timestamps; here we utilize a discrete approach and focus mostly on cascade sizes and publishers co-occurrences. Also, modelling of fluctuation scaling observed in online social communities is lastly a vivid topic of research (mostly with a random diffusion model [31], also with a time varying scale parameter to model a word occurrences fluctuations scaling in blogosphere [32]). Moreover, [33] and [34] already cover the topic of varying attractiveness of a message propagating in a complex network, however it is attributed to its producer not the message itself. To the best of our knowledge no other study addresses the temporal fluctuation in the news outlet network nor in the independent cascade model.

The rest of the document is structured as follows: first, we describe the dataset in detail and gives its basic statistics (Section 2); then we show that the media outlets activity follows a fluctuation scaling law, and report how scaling exponents depend on timescale and unit size (Section 3.1); next, we describe a reconstructed network of news outlets based on content correlations and basic features of such a system (Section 3.2); we finish the results section with showing how the independent cascade model can recover stylized facts observed in the data (Section 3.3). Finally, we draw conclusions, discuss, and summarize the work (Section 4). In three appendices, we consider basic statistical properties of investigated datasets (Appendix A), and introduce algorithms applied to compute fluctuation scaling exponents (Appendix B) and extract publishers network (Appendix C).

## 2. Data and its basic statistics

Event Registry platform monitors RSS feeds of over 25,000 news outlets for articles in 35 major languages from around the world, forms temporal clusters of similar articles (*events*), and extracts metadata about each event (such as involved entities, recognized Wikipedia-based concepts, location) [25]. A list of covered sources includes the biggest players in the news industry like online versions of tabloids (e.g. *dailymail.co.uk*) or national daily newspapers (e.g. *welt.de*), international news aggregators (e.g. *www.msn.com*), as well as numerous local or otherwise narrowly-focused news providers. A cluster of at least 5 similar articles, published within an interval between the newest and the oldest article of at most 4 days, forms an *event*. The dataset for the analysis consists of around 1,500,000 events (over 22,000,000 articles) published in 2016 and mentioning one of 11 subjectively chosen entities and concepts. We selected mostly keywords which were connected to major events of the year 2016. The first group were three major figures in the presidential campaign in the United States of America (the previous president *Barack Obama*, the Democratic Party’s presidential candidate – *Hillary Clinton*, and the Republican Party’s candidate –

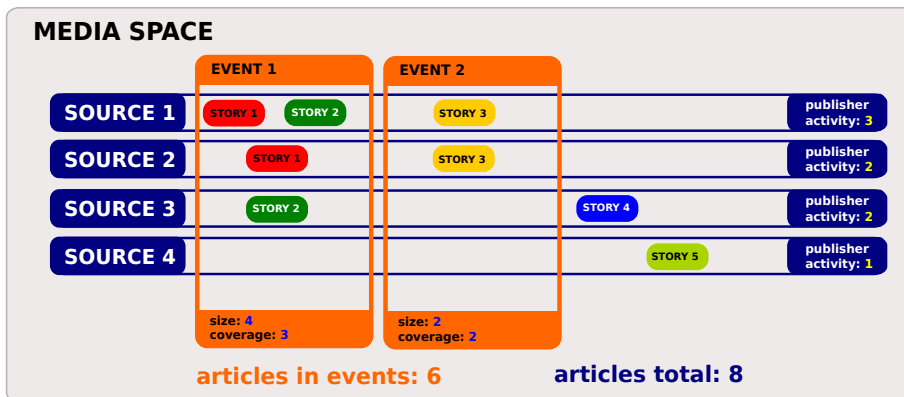


Figure 1: A visualization of the dataset structure and its considered features – event size (number of articles assigned to a given event), publisher activity (number of articles published by a given outlet), and event coverage (number of publishers having at least one article assigned to a given event). Story  $i$  is an article from an  $i$ -th cascade (see Methods) of very similar articles published by different sources. *Source 4* is an example of a publisher with no articles assigned to an event; other Sources are publishers in events. The Event Registry platform tracks only events consisting of at least 5 articles.

*Donald Trump*); the second group consisted of parties most actively involved in Brexit talks (*European Union, United Kingdom, Germany, France*). The *democracy* keyword alludes to both of the groups as those revolved around the voting process. *Association football* was selected because of its stable popularity in Europe and a major event in 2016 (UEFA European Championship). *Poland* and *Argentina* were selected as a kind of a baseline – both are countries not very visible in the international media, and devoid of globally impactful happenings in the year. In this Section we consider distributions of *publisher activities* (article number by publisher), and *event sizes* (article number by event) and *event coverages* (distinct publisher number by event). For a visual explanation – see Fig 1.

In Table 1, we show the number of articles and publishers (both total and assigned to events only) and number of articles assigned to an event for each concept. All columns in the table are strongly correlated with each other ( $r$ -Pearson coefficients  $> .9$  with  $p$ -values  $< 10^{-4}$ ), and the total number of articles is the best proxy to describe popularity of the keyword ( $r > .929$  with all the other columns) among the columns. The mean fraction of articles assigned to events is  $(43 \pm 2)\%$ , and the mean percent of publishers with at least one event-related article is  $(79 \pm 5)\%$  (both errors as a standard deviation in a given sample). Percents of articles in events and publishers in events have a relatively low sample-standard-deviation-to-mean ratio (below 0.06) comparing to other columns (above 0.12); it is not surprising, as the two mentioned values are normalized (thus intensive, opposing to the rest of the columns which contain extensive values). It also proves that, while a popularity varies among keywords, each keyword is covered in a similar way. Distributions of the aforementioned

concept	Articles	in events	Publishers	in events	events
Barack Obama	1,475,957	682,169 (46.2%)	12,067	9,941 (82.4%)	83,061
Hillary Clinton	1,302,358	586,636 (45.0%)	10,311	8,166 (79.2%)	54,095
Donald Trump	2,100,420	932,706 (44.4%)	11,700	9,506 (81.2%)	84,575
European Union	2,112,576	942,675 (44.6%)	11,989	9,740 (81.2%)	137,272
United Kingdom	3,432,419	1,497,603 (43.6%)	15,488	12,966 (83.7%)	265,734
Germany	3,706,909	1,546,498 (41.7%)	15,143	12,274 (81.1%)	278,536
France	3,247,340	1,397,352 (43.0%)	14,751	12,073 (81.8%)	227,404
Poland	511,151	207,339 (40.5%)	10,112	7,219 (71.4%)	42,850
Argentina	1,099,618	501,235 (45.6%)	10,357	7,788 (75.2%)	85,768
democracy	1,048,966	431,787 (41.1%)	11,520	8,974 (77.9%)	90,742
association football	2,321,356	970,251 (41.7%)	13,926	10,987 (78.9%)	184,988

Table 1: Basic properties of 11 examined concepts: numbers of articles, articles assigned to events, publishers, publishers which published at least one article assigned to an event, and events associated to each concept in the dataset. Percent values are calculated in relation to corresponding total numbers of articles or publishers.

values were all fat-tailed but we were unable to determine what function is the best description. For further considerations of the distributions fitting – see Appendix A.

### 3. Results

#### 3.1. Temporal fluctuation scaling in the dataset

The fluctuation scaling law (often called *Taylor’s law* after the famous L. R. Taylor’s paper [35]) is an empirical law observed in complex systems which consist of differently sized, otherwise similar, units. The law binds means and standard deviations of units’ activity (or, in general, a positive additive value characterizing each unit) in a form of a power law:  $\mu \sim \sigma^\alpha$ . The value of the exponent  $\alpha$  is known to indicate a degree of units’ synchronization in the observed system. It is common for  $\alpha$  to be in range from 0.5 (uncorrelated units) to 1.0 (perfect synchronization, e.g., strong external force) but also higher values have been reported [22] and theoretically approved [32]. Intermediate values are usually interpreted as a mixture of correlated and uncorrelated dynamics governing the system. There are two variants of the law. The ensemble fluctuation scaling can be observed if one groups units by a scale (*size-like*) parameter and consider means and variances calculated in each group. The temporal variant requires an observation of units’ activity which can be spotted over time. In this variant, the observation period is divided into time windows of size  $\Delta$ , then the activity is aggregated in each time window for each unit. The statistics  $\mu$  and  $\sigma$  are calculated for each unit separately over all time window [36]. In this study we focus on the latter variant and describe it in detail in Appendix B. Particularly, it is interesting for us to look at units’ activities aggregated in different time windows (i.e. at different timescales).

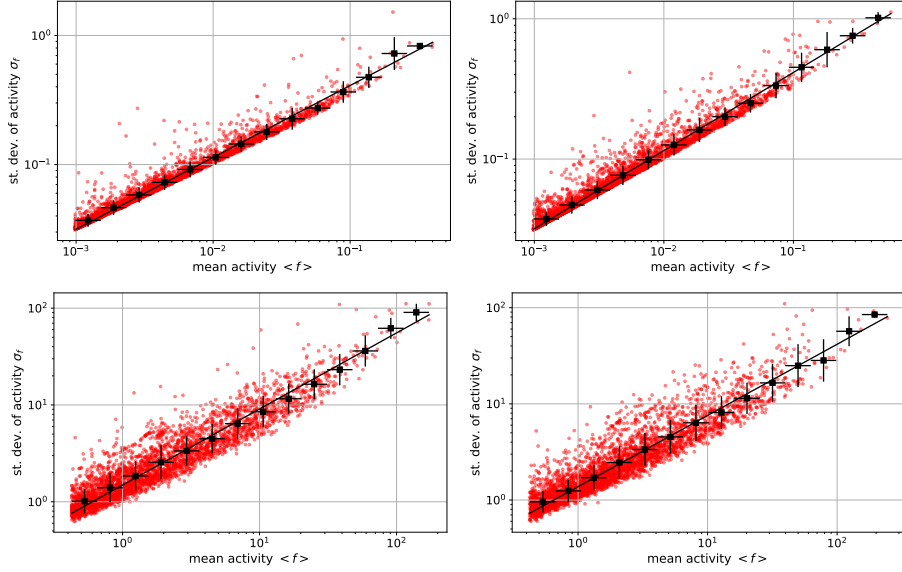


Figure 2: Temporal fluctuation scaling law of news outlets activity for time windows size  $\Delta \in \{10\text{min}, 3\text{days}\}$  for top and bottom, respectively. Articles about (left) *European Union* and (right) *association football*. X-axes – an average number of articles mentioning a given concept per fixed time window  $\Delta$  for a given publisher, Y-axes – a standard deviation of a number of the articles. Only publishers with more than 1 article per week on average were considered ( $N_{pubs} = 3349$  and  $N_{pubs} = 3974$ , respectively). Each red point represents one publisher, black points stand for a mean standard deviation in a given logarithmic bin, black line is a power law fit. Slopes are (left)  $\alpha(10\text{min}) = 0.557 \pm 0.022$ ,  $\alpha(3\text{days}) = 0.786 \pm 0.043$ , (right)  $\alpha(10\text{min}) = 0.557 \pm 0.013$ ,  $\alpha(3\text{days}) = 0.740 \pm 0.037$ .

We observed the temporal fluctuation scaling for all analyzed keyword and for all sizes of time windows. Figure 2 (left – *European Union*, right – *association football*) shows scatter plots of standard deviation  $\sigma$  versus mean  $\mu$  value of publisher activities (number of articles in a time window) for two selected time windows  $\Delta \in \{10\text{min}, 3\text{days}\}$ . The slopes of fitted lines are fluctuation scaling exponents  $\alpha(\Delta)$ . In each of the selected time window sizes, exponents are nearly the same for both concepts. For the longer time window, points are more scattered around the fitting line and the scaling exponents  $\alpha > 0.7$  (for the shorter one –  $\alpha \approx 0.55$ ).

The dependence of  $\alpha$  on  $\Delta$  for the aforementioned keywords is shown in Fig. 3 presenting three linear regimes with different slopes. A piecewise linear fit was applied to recover slopes in the regimes. Automatically detected breakpoints varied slightly for different keywords thus, for the sake of comparisons, we manually set them to 15min and 1day as these values were the most common in the automatic breakpoint detection. The exponent  $\alpha$  grows (nearly) monotonically with  $\Delta$ . For short time windows (up to  $\sim 15$  min),  $\alpha$  is close to 0.5 and growing with a slope  $\gamma_1 \approx 0.01$  per decade. For longer time windows the

growth is much faster ( $\gamma_2 \approx \gamma_3 \approx 0.09$ ). Regime slopes for all keywords in all regimes can be found in Fig. 4.

### 3.2. *Extracted publishers network*

In this section we will use Event Registry data to extract publishers network that will be further used in the next section to run a model of news cascades that reproduces the fluctuation scaling observed in the previous section. Nodes represent news outlets and weights of directed edges represent tie strengths. While not all news outlets can be observed and the exact propagation paths are impossible to follow, we hope to uncover meaningful connections between various publishers by processing reasonable amount of data related to contents of published articles. The full procedure of extracting the publishers network is described in Appendix C.

The network will be an environment for simulations of the independent cascade model (see the next section) thus we need to keep it reasonable size to make simulations faster. The resulting graph has 5,719 nodes and 1,329,030 edges. The vast majority of nodes was active only few times during the analyzed period, thus for readability and computational feasibility of simulations, we decided to prune the recovered network leaving only edges with  $u_{ij} > 0.5$  in the network recovered from the full dataset. The model was run on a giant component of the pruned graph (1,037 nodes and 4,150 edges).

A logarithmic binning of weighted degree distribution of the pruned graph (Fig. 5) shows that the in-degree distribution is wider than the out-degree distribution. Degrees of nodes in real networks are often power-law distributed but in the recovered graph there are nodes with a relatively high number of neighbors but not as much lowly connected nodes as one would expect. This is probably caused by the filtering of publishers with low activity.

A maximum spanning tree of the giant connected component of the pruned network is presented in Fig. 6 (only nodes with  $N_i > 5$ ). The MST was calculated for an undirected graph with weights set to  $\max(d_{ij}, d_{ji})$ . Geographical clustering of nodes for the major English-speaking countries (UK, USA, India, Australia, New Zealand) with addition of English versions of major local outlets (China, Africa) is clearly visible. While the visualization does not allow to determine direction of an edge, the most active connections reveal a few major information flow channels.

### 3.3. *Model*

The independent cascade model is an epidemic model run on a complex network. Nodes represents individuals (persons, social network users, news outlets) which can be in one of three possible states – susceptible, infected, or recovered. Directed edges represent probability of infection (information) transmission from infected to susceptible nodes. Infected nodes become recovered. The model has been recently extensively explored in the fields of social network analysis (to describe spreading of influence in social networks [37], to predict information diffusion probabilities [38]) and socio-physics (to model spreading

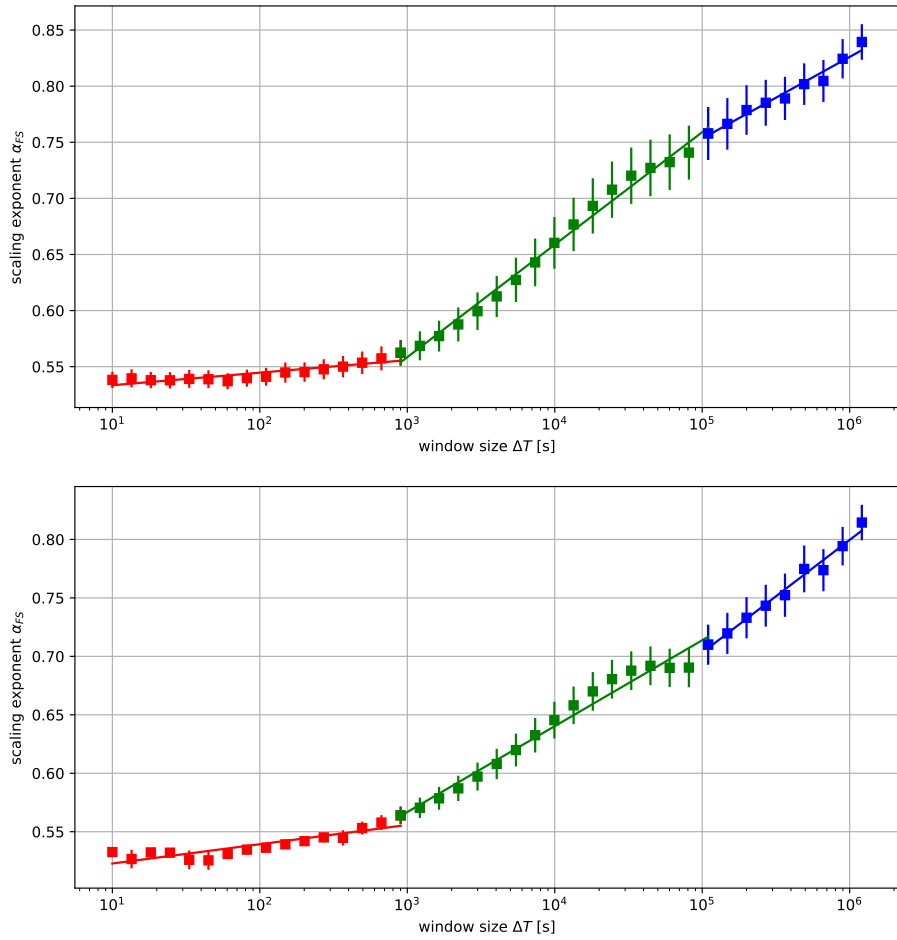


Figure 3: The temporal fluctuation scaling exponent:  $\alpha$  is nearly monotonically increasing with time window size  $\Delta$ ; character of the dependence is similar for all the analyzed concepts. (top) *European Union*, (bottom) *association football*. X-axis is in log scale. Lines are fit of logarithm function  $\alpha \sim \log \Delta$ . Breakpoints are set manually to  $\Delta = 15$  min and  $\Delta = 1$  day. Slopes are (sequentially red, green, and blue line): (left)  $\gamma_1 = 0.011 \pm 0.002$ ,  $\gamma_2 = 0.100 \pm 0.003$ ,  $\gamma_3 = 0.073 \pm 0.006$ , (right)  $\gamma_1 = 0.016 \pm 0.002$ ,  $\gamma_2 = 0.073 \pm 0.003$ ,  $\gamma_3 = 0.096 \pm 0.006$ .



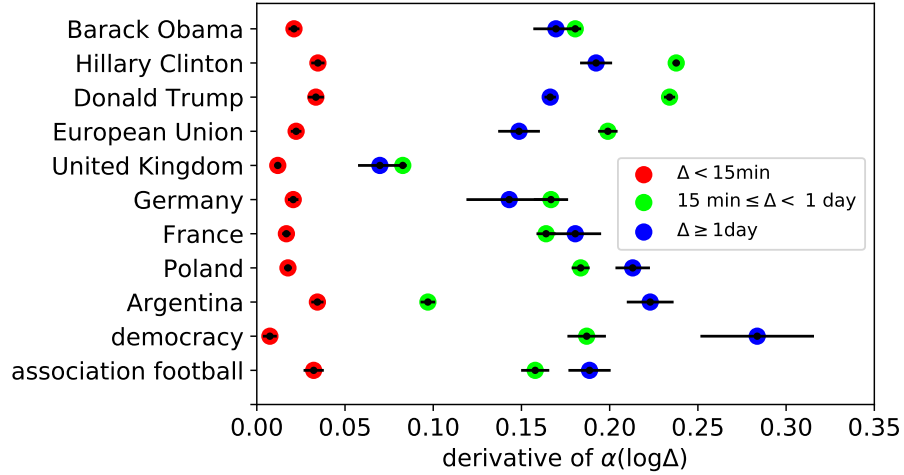


Figure 4: Slopes of temporal fluctuation scaling  $\alpha(\log \Delta)$  for different timescales and concepts. In shorter timescales ( $\Delta < 15\text{min}$ ) the slopes are similar for all the analyzed concepts. For longer timescales, the slopes are varied and, in few cases, separated for  $15\text{min} \leq \Delta \leq 1\text{day}$  and  $\Delta > 1\text{day}$ . Errors as a standard deviation of the slope coefficient.

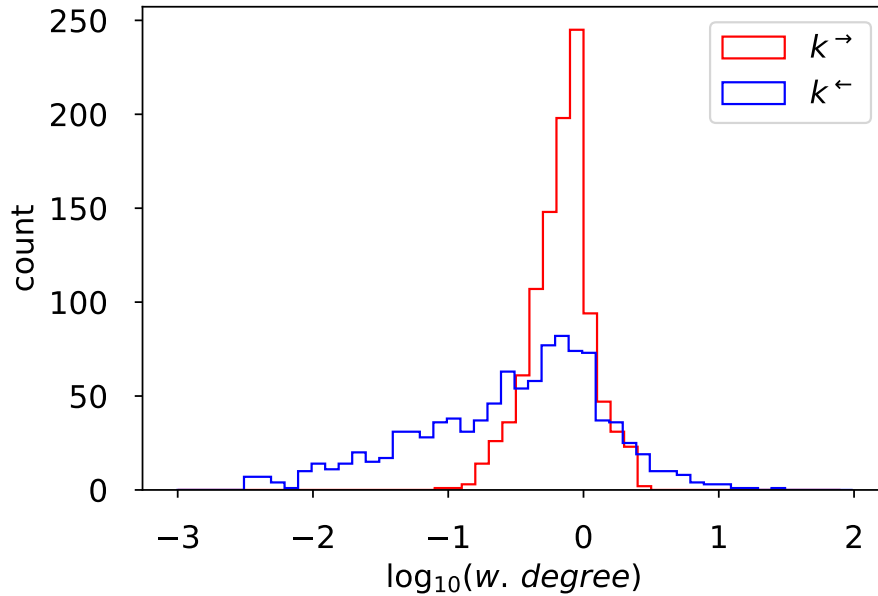


Figure 5: The histograms of weighted degree of nodes in the pruned publishers network do not follow a power law; in-degrees distribution is wider than the out-degrees distribution. Weighted in-  $k^{\leftarrow}$  and out-degrees  $k^{\rightarrow}$  are defined in Appendix C)). Logarithmic bins.

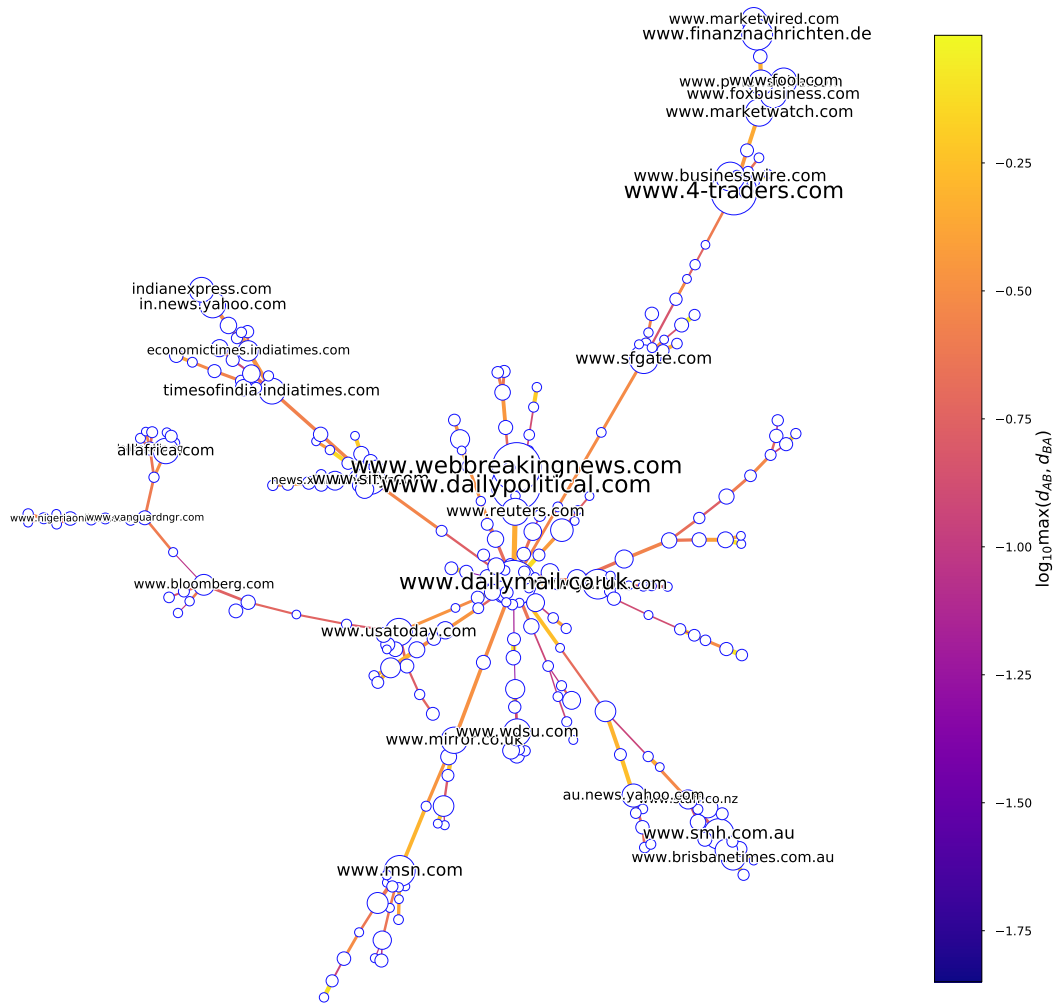


Figure 6: Extracted publishers network: the network consists of geographical and topical clusters; local hubs and important information flow channels are visible. Maximum spanning tree of the giant component of the extracted publishers network based on  $\max(d_{ij}, d_{ji})$ . Only publishers with  $N_i > 5$ ; only edges with  $\max(d_{ij}, d_{ji}) > 0.01$ , color of edge depends on logarithm of  $\max(d_{ij}, d_{ji})$ . Size of a node is proportional to its  $N_i$ . Width of an edge is proportional to its  $u_{ij}$ . Labels were shown for nodes with a degree in the MST over 2. The graph has 289 nodes and 288 edges.

of emotions [39] and reemergence of information diffusion [40]). Here we apply it to model the process of articles diffusion between news outlets. Heuristically, whenever an outlet publishes an original news item (or copies the article from an unobserved source), the competing outlets might decide to publish it as well which in turn might lead to another re-use of the article by their neighbors; moreover, if more than one neighbor of a node becomes infected then it is more likely the node will become infected which is a decent representation of a peer pressure between competitors.

We use the *independent\_cascade* add-on to the *networkx* [41] Python library written by Hung-Hsuan Chen. The outline of the algorithm using the SIR models-related terminology is as follows:

1. All nodes start *susceptible*,
2. Change status of one randomly chosen node (source) to *infected*,
3.  $t = 0$ ,
4. While(number of infected nodes  $> 0$ ):
  - (a)  $t = t + 1$ ,
  - (b) Each infected node A tries to infect each of its susceptible neighbors B with probability  $p_{AB}$ ,
  - (c) Each node infected in step  $t - 1$  becomes *recovered*.

The process was simulated on the network extracted from Event Registry data (*real*, see Methods), and two types of synthetic networks – a random graph (*ER* – Erdős-Rényi) and a Barabási-Albert network (*BA*). Sizes of artificial networks were set to be the same as in the giant component of the pruned network ( $N = 1,037$ ). The probability of connection between nodes in ER graph was set to be equal to the density of the component ( $p = \rho \approx 0.004$ ,  $\langle k \rangle = Np > 1$ ). In the BA graph, we assumed that the starting number of nodes  $d_0 = d$  and calculated  $d = \langle k \rangle / 2 \approx 4$ . Each edge in generated undirected graphs was changed to two directed edges - one in each of directions. We considered three above-mentioned networks with both homogeneous (*const*;  $d_{ij} = 0.02$ ) and heterogeneous edge weights (*shuffled*, drawn from the empirical distribution – for the three networks; *real*, calculated from the data – for the recovered network). This sums up to seven variants of a network topology and edge weights – *real real*, *real shuffled*, *real const*, *ER shuffled*, *ER const*, *BA shuffled*, *BA const*.

When we assumed that probabilities of link activation are the weights given by Eq. C.3, *i.e.*  $p_{AB} = d_{AB}$ , then simulated cascades were rather small (as should be expected for coverage of local and everyday news). To remove this discrepancy we further assumed that the news attractiveness might be captured by a real multiplicative factor  $h > 0$  (*hype*). The edge weights  $d_{AB}$  were multiplied by the selected hype:  $p_{AB} = \min(h d_{AB}, 1)$ . In Fig. 7, there are histograms of cascade sizes for the three networks with shuffled empirical weights with  $h \in \{1, 2, 5, 10, 20, 50\}$ .  $h = 2$  was enough to generate cascades of sizes comparable to half of the network size, and  $h = 50$  caused all nodes reachable from the source to be infected in every simulation (as expected –  $p_{AB} < 0.02$  were discarded). For higher  $h$ , more cascades exceeded the viral threshold and the expected size of the viral cascade was higher while its variance was lower.

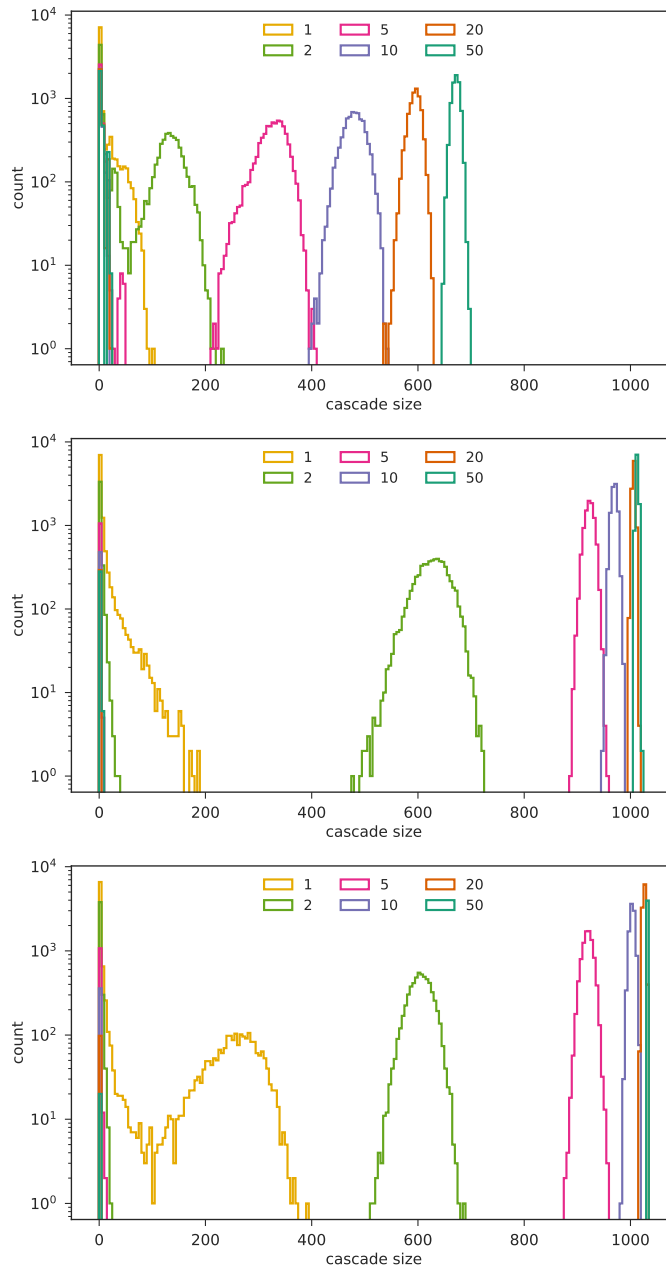


Figure 7: Cascade sizes distributions. It is visible that the hype parameter  $h$  controls expected size of cascade and chance to exceed viral threshold. Histograms of simulated cascade sizes for different values of multiplicative factor  $h$  on (top) the extracted network, (center) a random graph, (bottom) Barabasi-Albert graph. Y-axis is logarithmic. Simulations were performed 10,000 times in each setting.

We ran the independent cascade model on each of the selected networks and then performed the fluctuation scaling exponent calculations. We performed the simulations in batches of 10,000 cascades in a few variants of the hype parameter. Simulation in the first group of variants ( $C$ ) were conducted using the same  $h$  value for each cascade ( $Cx$  meaning  $h = x$ ; the simulations correspond to those presented in Fig. 7); in the second group of variants ( $P$ ) the hype parameter was selected at random from the power law distribution with the exponent  $\beta$  normed for hypes from 1 to 100 ( $Px$ :  $\beta = -x$ ); the third group ( $TP$ ) is similar to the second one but the simulations are performed in a sequence from the lowest to the highest values of  $h$  mimicking temporal correlations of the hype parameter; the last variant consisted of samples drawn from a uniform distribution in the range of 1–50 ( $uni$ ). The number of hype variants totals to 20 ( $C1, C2, C5, P1, P1.5, P2, P2.5, P3, P3.5, P4, P4.5, TP1, TP1.5, TP2, TP2.5, TP3, TP3.5, TP4, TP4.5, uni$ ).

The temporal fluctuation scaling was observed for all network variants but one ( $BA$  shuffled). Exemplary plots of the fluctuation scaling observed in *real real* network are shown in Fig. 8 ( $\Delta = 100$ , top –  $P4$ , bottom  $TP4$ ). Moreover, in the group  $TP$  for *real real*, *real shuffled*, and *ER shuffled* networks, we found that the fluctuation scaling exponent depends on the window size  $\Delta$  (Fig. 9). In all other cases, the scaling exponent was very close to 0.5 for all window sizes (*ER const*, *real const*), or the activity range was insufficient (less than one decade) to meaningfully recover  $\alpha$  (*BA const*, *BA shuffled*).

Interestingly, there is barely any difference between results for the recovered network and for the recovered network with shuffled weights; results are also similar for the *ER shuffled* but with a modest activity range (1.5 decade). The character of the dependence is similar to the one obtained for the real data in the previous sections. For networks with homogeneous edge weights, numerical simulations gave  $\alpha(\Delta) \approx 0.5$  which suggests that the existing diversity of interaction strengths between publishers is responsible for the observed dependence of the  $\alpha$  exponent on the length of the window size  $\Delta$  (see Fig.3).

We conclude the relationship between the fluctuation scaling exponent and the aggregation window size can be received from the independent cascade model run on a heterogeneous network with a slowly changing hype parameter controlling spreading rate.

#### 4. Discussion

In the paper, we present analyses of a dataset consisting of 22 million articles gathered by EventRegistry.org in 2016 containing at least one of 11 keywords published by over 10,000 news outlets from around the world.

First, we show long-tailed distributions found in the dataset. Distributions of number of articles published by different outlets can be described using the Weibull distribution; for distributions of number of publishers in different events, and event sizes the log-normal functions were the best fit.

Second, we consider the temporal fluctuation scaling of news outlets' activity around certain keywords. The result put the system on a long list of complex

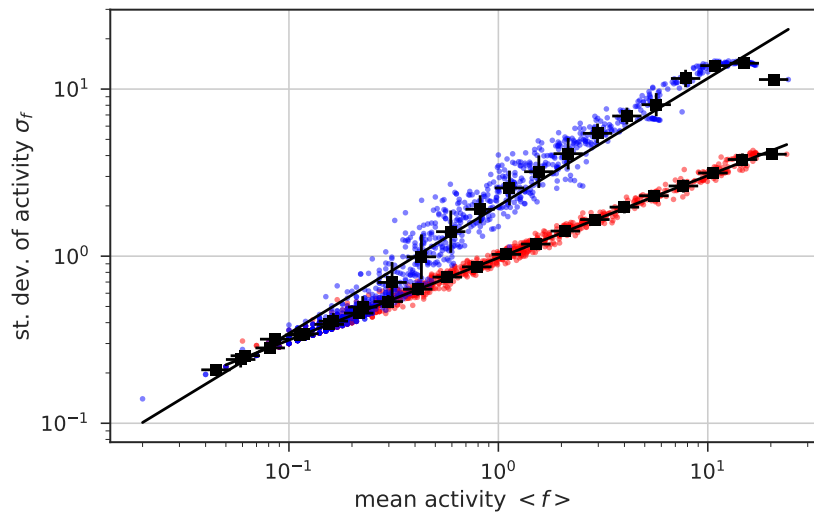


Figure 8: Evidence of fluctuation scaling in a model of publishers network. Plots for  $\Delta = 100$  and  $p(h) \sim h^{-4}$  ( $1 \leq h \leq 50$ ) (red points) without temporal hype correlations, (blue points) with temporal hype correlations. The recovered network with the recovered weights was an environment for both simulation batches. Each of simulation batches had 10,000 realizations. Slopes are (top)  $0.492 \pm 0.002$ , (bottom)  $0.763 \pm 0.025$ .

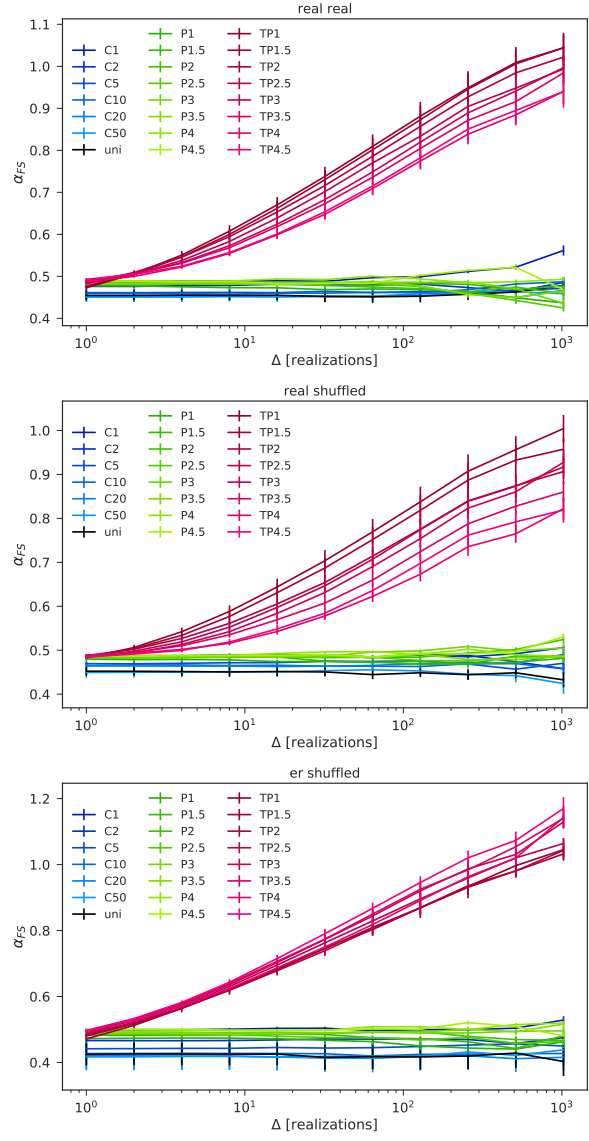


Figure 9: Temporal fluctuation scaling in model for different hype distributions ( $1 \leq h \leq 100$ ) in the recovered network (top), the recovered network with shuffled edge weights (middle), a random graph with empirical weights (bottom). Only models where edge weights were drawn from the empirical distribution and where simulations with similar  $h$  values were grouped (temporal clustering) indicate non-trivial fluctuation scaling exponents. Each color line stands for different hype distributions and orderings (see Section 3.3). Each of simulation batches had 10,000 realizations.

systems following the Taylor’s law [36]. The data follows nontrivial fluctuation scaling with three recognized regimes ( $\Delta < 15\text{min}$ ,  $15\text{min} < \Delta < 1\text{day}$ ,  $\Delta > 1\text{day}$ ). The result suggests that there are different dynamics governing different timescales – barely any correlations for short timescales, and a varying amount of synchronization for longer timescales.

Third, we uncover a network of correlations between news outlets content basing on co-occurrences in events and microclusters. The revealed network has interesting features – e.g. the in-degree distribution is wider than the out-degree distribution, geographical clustering, few strongly correlated groups of sources.

Fourth, we run the independent cascade model on the reverse engineered network and compare it to similar processes run on two synthetic networks (random graph and Barabasi-Albert). Although the independent cascade model leads to the temporal fluctuation scaling for nearly all cases, however to obtain nontrivial exponents observed in the data it is necessary to introduce a multiplicative parameter for a transmission chance (*hype*). Long-tailed event sizes distributions can be obtained using a long-tailed distribution of *hypes* as there is an expected cascade size for a given *hype* for a given network. Moreover, introducing grouping cascades with similar hype yields with an  $\alpha(\Delta)$  dependence similar as in the real data. We stress that the uncovered network gives a much better fit to the fluctuation scaling observed in the data as compared to the investigated synthetic networks.

The above analyses show a few interesting features of the dataset. Statistical inspection suggests the dynamics of news publishing is similar for each outlet depending mostly on a general activity of the outlet on a certain topic. The presence of long-tailed distribution seems to be a universal feature in human online communication channels [9, 42], or more generally speaking in complex systems [43]. The global news network follows the temporal fluctuation scaling law which is unsurprising as the system consists of spatially/temporally correlated units connected with overlapping communities [24]. Collective effects are stronger for longer timescales what corresponds to burstiness of media attention. Basing on our model the specific values of estimated scaling exponents are probably given by the structure of the underlying communication/mimicking network and the distribution of *hype* factor among stories. The proposed *hype* parameter might be interpreted as an external field coupled to observed outlets activity [22]. The results could be used to meaningfully estimate an impact of a given online story or an influence of a news outlet.

The presented model is surely not the only way to recover the fluctuation scaling similar to one observed in the data (in general there may be arbitrary many models indicating a given fluctuation scaling) but it shows a role of content attractiveness in information spreading and its fluctuation scaling and provides an interesting interpretation to the uncovered network. The method used to uncover a network of publishers might be an interesting tweak to existing methods of network recovery for cases when the original source or diffusion path is unclear. Our model focuses on propagation of news on a specific topic and does not take into account interests of news outlets. A model considering the whole news flow should definitely include eagerness of a given news outlet to



cover given topic (e.g. in a form of news’ topic vectors and outlets’ interests vectors [30]). Moreover, it would be interesting to broaden a range of analyzed time windows but it was impossible for the real data and very time-consuming for the model. It might be fruitful to consider temporal aspects of links between outlets or even treat the system as a coevolving network [44]. Also ego-networks of publishers and community structure of the network might be worth a closer look. An application of recent advancements in cross-lingual text comparisons (e.g. [45]) could lead to uncovering the global content correlation networks.

## Appendix A. Basic statistical properties of investigated datasets

We observed that logarithmically-binned density histograms of publisher activities, event sizes, and event coverages are long-tailed. To find the best fit in each case, we considered the following discrete positive ( $x \in \{0, 1, 2, \dots\}$ ) distributions provided by the *powerlaw* Python library [46]: power law with exponential cut-off  $f(x; \alpha, \kappa) \sim x^{-\alpha} e^{-\kappa x}$ , where  $\alpha, \kappa > 0$ ; positive log-normal  $f(x; \mu, \sigma) \sim x^{-1} \exp(-(\ln x - \mu)^2 / 2\sigma^2)$ , where  $\mu, \sigma > 0$ ; Weibull  $f(x; \beta, \lambda) \sim (x\lambda)^{\beta-1} \exp(-(\lambda x)^\beta)$ , where  $\beta, \lambda > 0$ .

Distribution parameters were calculated using corresponding maximum likelihood methods. To determine which of the distributions describes a given histogram most accurately, the fitted functions were compared pairwise using the log-likelihood ratios [47]. For any given pair, a likelihood of the data was calculated under each of competing distributions separately, then the Vuong’s log-likelihoods ratio test [48] was performed to determine which distribution was a better fit and whether the result was statistically significant.

For most concepts, the Weibull distribution was the better fit to histograms of publisher activities than the log-normal ( $p < 0.05$ ) and the truncated power-law ( $p < 0.05$ ) distributions. For the sake of comparability of the results among concepts, parameters for the Weibull distribution ( $\beta_A, \lambda_A$ ) were provided for all concepts. Figure A.10 shows histograms of publisher activity for keywords *European Union* and *association football*.

In case of event sizes, the statistical comparisons were inconclusive and we were unable to differentiate between the Weibull, log-normal, and truncated power law functions in all but two cases. The log-normal distribution ( $\mu_{EV}, \sigma_{EV}$ ) was chosen to be displayed in an aggregate table (Tab. A.2). Figure A.11 presents histograms of sizes of event about *European Union* and *associated football* in terms of article count (event size).

For the majority of the analyzed concepts the log-normal distribution was the best fit to the empirical distribution of event coverages (with  $p < 0.05$  for 7 out of 11 concepts) thus  $\mu_{EC}$  and  $\sigma_{EC}$  are provided for each concept for a comparison. Histograms in Fig. A.12 present distributions of sizes of events about *European Union* and *associated football* in terms of involved publishers.

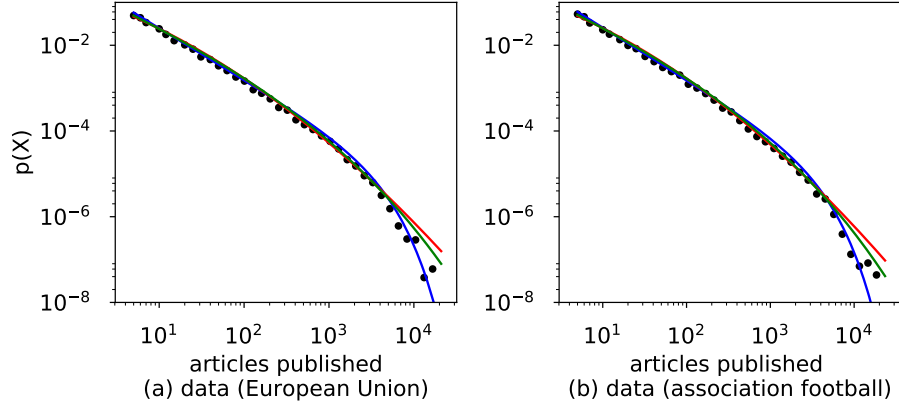


Figure A.10: Normed histograms of **publisher activities** for concepts (left) *European Union* and (right) *association football* published by each source. X-axis – a number of articles with a given concept published by a source, Y-axis – a normalized count. Data grouped in logarithmic bins, color lines are various types of fitted heavy-tailed distributions (blue – truncated power law, red – log-normal, green – Weibull). The Weibull distribution turned out to be the best fit across majority of analyzed keywords. Fit parameters are in Tab. A.2

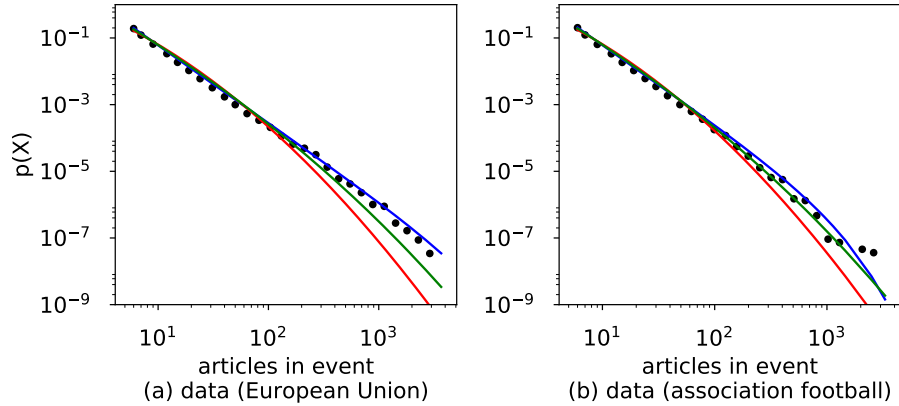


Figure A.11: Normed histograms of **sizes of events** containing concepts (left) *European Union* or (right) *association football*. X-axis – a number of articles with a given concept assigned to an event, Y-axis – a normalized count. Data grouped in logarithmic bins, color lines are various types of fitted heavy-tailed distributions (blue – truncated power law, red – log-normal, green – Weibull). The log-normal distribution was selected the best fit across majority of analyzed keywords. Fit parameters are in Tab. A.2

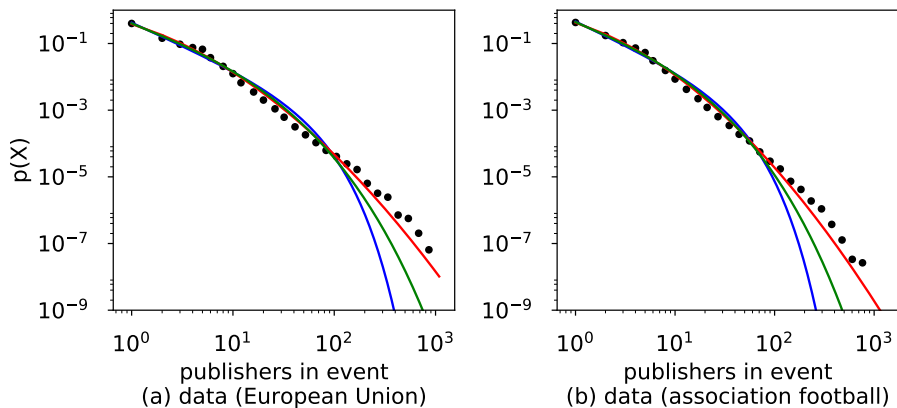


Figure A.12: Normed histograms of **coverage of events** containing concepts *European Union* (left) and *association football* (right). X-axis – a number of sources which published an article assigned to an event, Y-axis – a normalized count. Data grouped in logarithmic bins, color lines are various types of fitted heavy-tailed distributions (blue – truncated power law, red – log-normal, green – Weibull). The log-normal distribution was selected the best fit across majority of analyzed keywords. Fit parameters are in Tab. A.2

concept	$\beta_A$	$\lambda_A$	$\mu_{EV}$	$\sigma_{EV}$	$\mu_{EC}$	$\sigma_{EC}$
Barack Obama	0.25	0.35	$1.9 \times 10^{-8}$	1.70	$3.4 \times 10^{-8}$	1.66
Hillary Clinton	0.23	0.67	$4.1 \times 10^{-7}$	1.74	$4.7 \times 10^{-8}$	1.77
Donald Trump	0.24	0.33	$1.0 \times 10^{-8}$	1.72	0.0636	1.79
European Union	0.24	0.28	$1.4 \times 10^{-7}$	1.50	0.2113	1.49
United Kingdom	0.28	0.08	$7.0 \times 10^{-8}$	1.48	$1.8 \times 10^{-9}$	1.46
Germany	0.26	0.15	$1.0 \times 10^{-7}$	1.48	0.0083	1.46
France	0.27	0.11	$2.2 \times 10^{-8}$	1.50	0.0762	1.49
Argentina	0.21	1.57	$6.7 \times 10^{-8}$	1.45	0.3077	1.41
Poland	0.25	1.00	$5.9 \times 10^{-7}$	1.45	$1.5 \times 10^{-7}$	1.38
democracy	0.27	0.24	$4.2 \times 10^{-7}$	1.40	0.0359	1.40
association football	0.27	0.15	$1.8 \times 10^{-7}$	1.44	0.1149	1.37

Table A.2: Fitted parameters of distributions.  $\beta_A, \lambda_A$  – Weibull fit parameters for distribution of publisher activities,  $\mu_{EV}, \sigma_{EV}$  – log-normal fit parameters for distribution of event sizes,  $\mu_{EC}, \sigma_{EC}$  – log-normal fit parameters for distribution of events coverage (number of unique publishers).

## Appendix B. Temporal fluctuation scaling

The temporal fluctuation scaling has been applied to our data as follows. Let  $f_{i,t}$  be a positive variable describing an additive measure of an activity of the object  $i$  at time moment  $t$ . Examples of such activities can be a number of data packages coming to a router, emails sent by a person, or articles published by an outlet. Let the total number of elements in time series of this activity be  $T$ , i.e.  $t = 1, 2, 3, \dots, T$  (further we will assume that  $T$  is the same for all units  $i$ ). Let us further divide the series into  $Q$  windows of size  $\Delta$ , i.e.,  $Q\Delta = T$ . The quantity  $f_i^{(q,\Delta)}$  stands for a cumulative value of the variable  $f_i$  in a window of size  $\Delta$  ( $q = 1, 2, 3, \dots, Q$  is the window's label) and  $(\sigma_i^\Delta)^2$  is the variance of this cumulative variable in the whole data series. Then we have

$$(\sigma_i^\Delta)^2 = \langle [f_i^{(q,\Delta)}]^2 \rangle - \langle [f_i^{(q,\Delta)}] \rangle^2 \quad (\text{B.1})$$

Here

$$\langle [f_i^{(q,\Delta)}] \rangle = Q^{-1} \sum_{q=1}^Q \sum_{t=(q-1)\Delta+1}^{q\Delta} f_{i,t} = \Delta \frac{\sum_{t=1}^T f_{i,t}}{T} \quad (\text{B.2})$$

and

$$\langle [f_i^{(q,\Delta)}]^2 \rangle = 1/Q \sum_{q=1}^Q \left( \sum_{t=(q-1)\Delta+1}^{q\Delta} f_{i,t} \right)^2 \quad (\text{B.3})$$

It was observed for router activity and email traffic (but also stock markets, river flows, or printing activity) [36] that:

$$\sigma_i^\Delta \propto \langle [f_i^{(q,\Delta)}] \rangle^{\alpha(\Delta)}. \quad (\text{B.4})$$

In practice, data loosely follows the Taylor scaling law thus in order to estimate the exponent  $\alpha(\Delta)$  the following procedure was applied:

1. Calculate mean  $\log\langle f \rangle$  for data equally binned by  $\log \sigma_i$ ,
2. Perform least squares fit to the binned data.

The exponent  $\alpha$  usually depends on  $\Delta$  and few linear regimes can be observed – the systems observed in [36] followed two, and news outlets activity as described in this study consistently followed three. We used the piecewise linear fit Python library *pwlfit* which is the C. Jekel's implementation of the Least-squares Fit of a Continuous Piecewise Linear Function [49].

To guarantee sensible statistics for all analyzed units (publishers), we discarded those which had on average less than one article mentioning a given keyword per week in the analyzed period (thus the activity threshold is equal to 52). Units with mean activity below such a threshold also follow the fluctuation scaling law but with the trivial exponent  $\alpha = 0.5$ ; the effect is caused by a relative sparsity of the signals [50, 32].

## Appendix C. Extracting publishers network

A common approach to uncover the underlying propagation network would be to use information about publication time [29]. Because of the incomplete knowledge about all sources of articles and not fully reliable timestamps, in our case it is hard to determine the original source of a given piece of content [51]. We decided to use a co-occurrence fraction counting method known from the field of scientometrics [52]. Thus we calculated similarities between articles in each event to find clusters of highly similar articles (*cascades*) to track which news outlets frequently co-occur in the cascades. Event Registry clustering functionality reduced the required number of pairwise comparisons by a few orders of magnitude.

To extract cascades from the data, the following procedure was applied to each event from the given period:

1. download articles,
2. transform each article to a vector of 3-gram occurrences with TF-IDF weighting (trained on a set of 10,000 randomly selected articles from a week preceding the publication date),
3. calculate a cosine similarity matrix and use it as a distance matrix for a single-linkage hierarchical clustering of the articles,
4. obtain clusters at the threshold value set to 0.25 (it should guarantee reliable results – see [27]),
5. save a list of unique publishers in each cluster.

The procedure was implemented using the *scikit-learn* Python library [53]. We decided to use a list of unique publishers because many similar articles from the same source might be caused by the crawler malfunction (e.g. incorrectly obtained article body) or a few updates of the same text. Event Registry system uses filtered stems, concepts, and other metadata; here we used 3-grams to focus on correlations of the content (not only topic) of the articles.

The publishers network was modeled using all 14,851 events with articles in English language published between 01-08.05.2017. In the dataset, we found 22,525 cascades with at least two involved publishers (the total number of cascades was 168,725).

Let  $C$  be a number of observed cascades,  $P$  – total number of publishers contributing. Following [52], we construct an occurrence counting matrix  $A = [a_{pc}]$  with  $P$  rows and  $C$  columns. The matrix  $A$  can be also seen as an adjacency matrix of a bipartite graph where the two types of nodes are publishers and cascades (like papers and authors [54], people and social media platforms [55], companies and economic sectors [56]). Each element of matrix  $A$  is defined as:

$$a_{pc} = \begin{cases} 1, & \text{if the } p\text{-th publisher has an article in the } c\text{-th cascade.} \\ 0, & \text{otherwise.} \end{cases} \quad (\text{C.1})$$

The above definition allows us to retrieve a fractional counting cascade co-occurrence  $P \times P$  symmetric matrix  $U = [u_{ij}]$  as:

$$u_{ij} = \sum_{c=1}^C \frac{a_{ic}a_{jc}}{n_c^2} \quad (\text{C.2})$$

where  $n_c$  – size of the  $c$ -th cascade. Moreover, the diagonal elements  $u_{ii} = N_i$  are fractional occurrence counts for the  $i$ -th publisher.

We define an asymmetric matrix  $D = [d_{ij}]$  representing a weighted adjacency matrix of the directed publishers network:

$$d_{ij} = u_{ij}/N_i \quad (\text{C.3})$$

Trivially,  $(\forall_{i,j})(N_i \geq u_{ij})$  as each publisher co-occurred with itself in all cascades it was involved. This means  $d_{ij} \leq 1$  and it allows us to use the matrix  $D$  as an input for the independent cascade model where  $d_{ij}$  will be further used to calculate a probability of an activation of a directed edge from  $i$ -th to  $j$ -th publisher assuming that the  $i$ -th publisher was infected in the previous model step.

The weighted out-degree of the  $i$ -th node is:

$$k_i^{\rightarrow} = \sum_{p=1, p \neq i}^P d_{ip} \quad (\text{C.4})$$

and it is equal to an expected number of edges activated by the  $i$ -th node at time  $t + 1$  if it was infected at time  $t$ . On the other hand the weighted in-degree of the  $i$ -th node is:

$$k_i^{\leftarrow} = \sum_{p=1, p \neq i}^P d_{pi} \quad (\text{C.5})$$

and it is equal to an average number of times it would be infected at time  $t + 1$  if its nearest neighborhood is completely infected at time  $t$ .

## Acknowledgements

This research has received funding as *RENOIR* Project from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 691152, by Ministry of Science and Higher Education (Poland), grant Nos. W34/H2020/2016, 329025/PnH/2016, and by National Science Centre, Poland Grant No. 2015/19/B/ST6/02612. J.A.H. was partially supported by the Russian Science Foundation, Agreement #17-71-30029 with co-financing of Bank Saint Petersburg. This research was also supported in part by PLGrid Infrastructure.

## Data Availability

The data that support the findings of this study are available from Event Registry but restrictions apply to the availability of these data, which were used under license for the current study, and so are not publicly available. Data are however available from the authors upon reasonable request and with permission of Event Registry.

## References

### References

- [1] R. Conte, N. Gilbert, G. Bonelli, C. Cioffi-Revilla, G. Deffuant, J. Kertesz, V. Loreto, S. Moat, J. P. Nadal, A. Sanchez, A. Nowak, A. Flache, M. San Miguel, D. Helbing, Manifesto of computational social science, *The European Physical Journal Special Topics* 214 (2012) 325–346.
- [2] C. Castellano, S. Fortunato, V. Loreto, Statistical physics of social dynamics, *Reviews of Modern Physics* 81 (2009) 591–646.
- [3] J. Kwapien, S. Drozd, Physical approach to complex systems, *Physics Reports* 515 (2012) 115 – 226. Physical approach to complex systems.
- [4] D. Helbing, D. Brockmann, T. Chadeaux, K. Donnay, U. Blanke, O. Woolley-Meza, M. Moussaid, A. Johansson, J. Krause, S. Schutte, M. Perc, Saving human lives: What complexity science and information systems can contribute, *Journal of Statistical Physics* 158 (2014) 735–781.
- [5] M. R. D’Orsogna, M. Perc, Statistical physics of crime: A review, *Physics of Life Reviews* 12 (2015) 1–21.
- [6] Z. Wang, C. T. Bauch, S. Bhattacharyya, A. d’Onofrio, P. Manfredi, M. Perc, N. Perra, M. Salathé, D. Zhao, Statistical physics of vaccination, *Physics Reports* 664 (2016) 1–113.
- [7] M. Perc, J. J. Jordan, D. G. Rand, Z. Wang, S. Boccaletti, A. Szolnoki, Statistical physics of human cooperation, *Physics Reports* 687 (2017) 1–51.
- [8] A. Chmiel, J. Sienkiewicz, M. Thelwall, G. Paltoglou, K. Buckley, A. Kappas, J. A. Holyst, Collective emotions online and their influence on community life, *PLOS ONE* 6 (2011) 1–8.
- [9] A. Garas, D. Garcia, M. Skowron, F. Schweitzer, Emotional persistence in online chatting communities, *Scientific Reports* 2 (2012) 402.
- [10] J. A. Holyst (Ed.), *Cyberemotions*, Springer International Publishing, International, 2017. doi:10.1007/978-3-319-43639-5.

- [11] P. Sobkowicz, M. Kaschesky, G. Bouchard, Opinion mining in social media: Modeling, simulating, and forecasting political opinions in the web, *Government Information Quarterly* 29 (2012) 470 – 479. *Social Media in Government - Selections from the 12th Annual International Conference on Digital Government Research* (dg.o2011).
- [12] T. Kuhn, M. Perc, D. Helbing, Inheritance patterns in citation networks reveal scientific memes, *Physical Review X* 4 (2014).
- [13] M. V. Tomasello, G. Vaccario, F. Schweitzer, Data-driven modeling of collaboration networks: a cross-domain analysis, *EPJ Data Science* 6 (2017) 22.
- [14] A. Patania, G. Petri, F. Vaccarino, The shape of collaborations, *EPJ Data Science* 6 (2017) 18.
- [15] J. Sienkiewicz, K. Soja, J. A. Hołyst, P. M. A. Sloot, Categorical and geographical separation in science, *Scientific Reports* 8 (2018) 8253.
- [16] A. M. Petersen, J. N. Tenenbaum, S. Havlin, H. E. Stanley, M. Perc, Languages cool as they expand: Allometric scaling and the decreasing need for new words, *Scientific Reports* 2 (2012).
- [17] M. Gomez-Rodriguez, J. Leskovec, A. Krause, Inferring networks of diffusion and influence, *ACM Transactions on Knowledge Discovery from Data* 5 (2012) 1–37.
- [18] M. G. Rodriguez, J. Leskovec, B. Schölkopf, Structure and dynamics of information pathways in online media, in: *Proceedings of the sixth ACM international conference on Web search and data mining - WSDM '13*, ACM Press, Rome, Italy, 2013, pp. 23–32. doi:10.1145/2433396.2433402.
- [19] N. Liu, H. An, X. Gao, H. Li, X. Hao, Breaking news dissemination in the media via propagation behavior based on complex network theory, *Physica A: Statistical Mechanics and its Applications* 453 (2016) 44 – 54.
- [20] X. He, Y.-R. Lin, Measuring and monitoring collective attention during shocking events, *EPJ Data Science* 6 (2017) 30.
- [21] M. Jalili, M. Perc, Information cascades in complex networks, *Journal of Complex Networks* 5 (2017) 665–693.
- [22] A. Fronczak, P. Fronczak, Origins of Taylor’s power law for fluctuation scaling in complex systems, *Phys. Rev. E* 81 (2010) 066112.
- [23] J. Chołoniowski, A. Chmiel, J. Sienkiewicz, J. A. Hołyst, D. Küster, A. Kappas, Temporal Taylor’s scaling of facial electromyography and electrodermal activity in the course of emotional stimulation, *Chaos, Solitons & Fractals* 90 (2016) 91–100.



- [24] G. Petri, P. Expert, H. J. Jensen, J. W. Polak, Entangled communities and spatial synchronization lead to criticality in urban traffic, *Scientific Reports* 3 (2013) 1798.
- [25] G. Leban, B. Fortuna, J. Brank, M. Grobelnik, Event Registry: Learning about world events from news, in: *Proceedings of the 23rd International Conference on World Wide Web, WWW '14 Companion*, 2014, pp. 107–110. doi:10.1145/2567948.2577024.
- [26] C. Jacobi, W. van Atteveldt, K. Welbers, Quantitative analysis of large amounts of journalistic texts using topic modelling, *Digital Journalism* 4 (2016) 89–106.
- [27] J. Chołojewski, G. Leban, A. Rehar, S. Maček, Information flow between news articles: Slovene media case study, in: *Proceedings of the 19th International Multiconference "Information Society"*, vol. D, 2016, pp. 13–16.
- [28] N. Barbieri, F. Bonchi, G. Manco, Cascade-based community detection, in: *Proceedings of the sixth ACM international conference on Web search and data mining - WSDM 13*, ACM Press, Rome, Italy, 2013, pp. 33–42. doi:10.1145/2433396.2433403.
- [29] J. Pouget-Abadie, T. Horel, Inferring graphs from cascades: A sparse recovery framework, in: *Proceedings of the 24th International Conference on World Wide Web - WWW 15 Companion*, ACM Press, Florence, Italy, 2015, pp. 625–626. doi:10.1145/2740908.2744107.
- [30] M. Yu, V. Gupta, M. Kolar, An influence-receptivity model for topic based information cascades, in: *2017 IEEE International Conference on Data Mining (ICDM)*, IEEE, New Orleans, LA, USA, 2017, pp. 1141–1146. doi:10.1109/icdm.2017.152.
- [31] Y. Sano, H. Takayasu, M. Takayasu, Fluctuation scaling in online social media, in: *2015 International Conference on Noise and Fluctuations (ICNF)*, IEEE, Xian, China, 2015, pp. 1–5. doi:10.1109/icnf.2015.7288568.
- [32] H. Watanabe, Y. Sano, H. Takayasu, M. Takayasu, Statistical properties of fluctuations of time series representing appearances of words in nationwide blog data and their applications: An example of modeling fluctuation scalings of nonstationary time series, *Phys. Rev. E* 94 (2016) 052317.
- [33] L. Alessandretti, K. Sun, A. Baronchelli, N. Perra, Random walks on activity-driven networks with attractiveness, *Physical Review E* 95 (2017) 052318.
- [34] I. Pozzana, K. Sun, N. Perra, Epidemic spreading on activity-driven networks with attractiveness, *Physical Review E* 96 (2017) 042310.
- [35] L. R. Taylor, Aggregation, variance and the mean, *Nature* 189 (1961) 732–735.

- [36] Z. Eisler, I. Bartos, J. Kertész, Fluctuation scaling in complex systems: Taylor’s law and beyond, *Adv. Phys.* 57 (2008) 89.
- [37] Z. Chen, K. Taylor, Modeling the spread of influence for independent cascade diffusion process in social networks, in: 2017 IEEE 37th International Conference on Distributed Computing Systems Workshops (ICDCSW), IEEE, Atlanta, GA, USA, 2017, pp. 151–156. doi:10.1109/icdcs.2017.28.
- [38] D. Varshney, S. Kumar, V. Gupta, Predicting information diffusion probabilities in social networks: A bayesian networks based approach, *Knowledge-Based Systems* 133 (2017) 66–76.
- [39] X. Xiong, Y. Li, S. Qiao, N. Han, Y. Wu, J. Peng, B. Li, An emotional contagion model for heterogeneous social media with multiple behaviors, *Physica A: Statistical Mechanics and its Applications* 490 (2018) 185–202.
- [40] D. Yang, X. Liao, H. Shen, X. Cheng, G. Chen, Modeling the reemergence of information diffusion in social network, *Physica A: Statistical Mechanics and its Applications* 490 (2018) 1493–1500.
- [41] A. A. Hagberg, D. A. Schult, P. J. Swart, Exploring network structure, dynamics, and function using networkx, in: G. Varoquaux, T. Vaught, J. Millman (Eds.), *Proceedings of the 7th Python in Science Conference*, Pasadena, CA USA, 2008, pp. 11 – 15.
- [42] J. Sienkiewicz, M. Skowron, G. Paltoglou, J. A. Hołyst, Entropy-growth-based model of emotionally charged online dialogues, *Advances in Complex Systems* 16 (2013) 1350026.
- [43] D. Sornette, Probability distributions in complex systems, in: R. A. Meyers (Ed.), *Encyclopedia of Complexity and Systems Science*, Springer-Verlag, New York, 2009, pp. 7009–7024. URL: <http://arxiv.org/abs/0707.2194>.
- [44] J. Toruniewska, K. Kułakowski, K. Suchecki, J. A. Hołyst, Coupling of link- and node-ordering in the coevolving voter model, *Physical Review E* 96 (2017) 042306.
- [45] J. Rupnik, A. Muhic, G. Leban, P. Skraba, B. Fortuna, M. Grobelnik, News across languages - cross-lingual document similarity and event tracking, *Journal of Artificial Intelligence Research* 55 (2016) 283–316.
- [46] J. Alstott, E. Bullmore, D. Plenz, powerlaw: A Python package for analysis of heavy-tailed distributions, *PLoS ONE* 9 (2014) 1–11.
- [47] A. Clauset, C. R. Shalizi, M. E. J. Newman, Power-law distributions in empirical data, *SIAM Review* 51 (2009) 661–703.

- [48] Q. H. Vuong, Likelihood ratio tests for model selection and non-nested hypotheses, *Econometrica* 57 (1989) 307.
- [49] N. Golovchenko, Least-squares fit of a continuous piecewise linear function, (2004). <http://www.golovchenko.org/docs/ContinuousPiecewiseLinearFit.pdf>.
- [50] S. Meloni, J. Gómez-Gardeñes, V. Latora, Y. Moreno, Scaling breakdown in flow fluctuations on complex networks, *Phys. Rev. Lett.* 100 (2008) 208701.
- [51] M. Kivelä, R. K. Pan, K. Kaski, J. Kertész, J. Saramäki, M. Karsai, Multiscale analysis of spreading in a large communication network, *Journal of Statistical Mechanics: Theory and Experiment* 2012 (2012) P03005.
- [52] L. Leydesdorff, H. W. Park, Full and fractional counting in bibliometric networks, *Journal of Informetrics* 11 (2017) 117–120.
- [53] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* 12 (2011) 2825–2830.
- [54] M. E. J. Newman, Scientific collaboration networks. I. Network construction and fundamental results, *Physical Review E* 64 (2001) 016131.
- [55] M. Mitrović, B. Tadić, Dynamics of bloggers’ communities: Bipartite networks from empirical data and agent-based modeling, *Physica A: Statistical Mechanics and its Applications* 391 (2012) 5264–5278.
- [56] A. M. Chmiel, J. Sienkiewicz, K. Suchecki, J. A. Hołyst, Networks of companies and branches in Poland, *Physica A: Statistical Mechanics and its Applications* 383 (2007) 134–138.