

The dynamics of information-driven coordination phenomena: a transfer entropy analysis

Short title: Dynamics of coordination phenomena

Javier Borge-Holthoefer,^{1*} Nicola Perra,^{2*} Bruno Gonçalves,^{3†} Sandra González-Bailón,⁴ Alex Arenas,^{5*} Yamir Moreno,^{6,7,8*} Alessandro Vespignani^{2,8,9*}

¹ Qatar Computing Research Institute, HBKU, Doha, Qatar

² Laboratory for the Modeling of Biological and Socio-technical Systems, Northeastern University, Boston 02115, USA

³ Aix Marseille Université, Université de Toulon, CNRS, CPT, UMR 7332, 13288 Marseille, France

⁴ Annenberg School for Communication, University of Pennsylvania, Philadelphia 19104, USA

⁵ Departament d'Enginyeria Informàtica i Matemàtiques, Universitat Rovira i Virgili, 43007 Tarragona, Spain

⁶ Department of Theoretical Physics, Faculty of Sciences, University of Zaragoza, Zaragoza 50009, Spain

⁷ Institute for Biocomputation and Physics of Complex Systems (BIFI), University of Zaragoza, 50018 Zaragoza, Spain Institute for Scientific Interchange (ISI), Torino, Italy

⁸ Institute for Scientific Interchange (ISI), Torino, Italy

⁹ Institute for Quantitative Social Sciences at Harvard University, Cambridge MA 02138, USA

* To whom correspondence should be addressed: jborge@qf.org.qa, n.perra@neu.edu, alexandre.arenas@urv.cat, yamir.moreno@gmail.com, a.vespignani@neu.edu

Abstract

Data from social media are providing unprecedented opportunities to investigate the processes that rule the dynamics of collective social phenomena. Here, we consider an information theoretical approach to define and measure the temporal and structural signatures typical of collective social events as they arise and gain prominence. We use the symbolic transfer entropy analysis of micro-blogging time series to extract directed networks of influence among geolocalized sub-units in social systems. This methodology captures the emergence of system-level dynamics close to the onset of socially relevant collective phenomena. The framework is validated against a detailed empirical analysis of five case studies. In particular, we identify a change in the characteristic time-scale of the information transfer that flags the onset of information-driven collective phenomena. Furthermore, our approach identifies an order-disorder transition in the directed network of influence between social sub-units. In the absence of a clear exogenous driving, social collective phenomena can be represented as endogenously-driven structural transitions of the information transfer network. This study provides results that can help define models and predictive algorithms for the analysis of societal events based on open source data.

Introduction

A vivid scientific and popular media debate has recently centered on the role that social networking tools play in coordinating collective phenomena. Examples include street protests, civil unrests, consensus formation, or the emergence of electoral preferences. A flurry of studies have analyzed the correlation of

search engine queries, microblogging posts and other open data sources with the incidence of infectious disease [1-4], box office returns [5], stock market behavior [6,7], election outcomes [8-9], popular votes results [10], crowd sizes [11], and social unrest [12,13]. Many other studies, however, have also pointed out the challenges big data presents and the likely methodological pitfalls that might result from their analysis [14-19]. This prior work suggests that more research is needed to develop methods for exploiting the value of social media data while overcoming their limitations.

Here, we use micro-blogging data to extract networks of causal influence among different geographical sub-units before, during, and after collective social phenomena. In order to ground our work on empirical data, we analyze five datasets that track Twitter communications around five well-known social events: the release of a Hollywood blockbuster movie; two massive political protests; the discovery of the Higgs boson; and the acquisition of Motorola by Google. We selected these case studies because they represent different points in a theoretical continuum that separates two types of collective phenomena: those that can be represented as an endogenously-driven exchange of information; and those that respond more clearly to factors that are exogenous to the system. In our context, these phenomena refer to dynamics of information exchange through social media: in some cases, discussions evolve organically, building up momentum up to the point where the exchange of information is generalized; in some other cases, however, the discussions emerge suddenly as a reaction to some unexpected external event [20]. The Motorola-Google case study corresponds to the exogenous type, providing a counterexample and intuitive baseline test for the rest.

For each case study we adopt the transfer entropy approach to define an effective social connectivity at the macro-scale, and study the coordinated activation of localized populations. We address two foundational problems: first, the identification of the characteristic time-scale of social events as they develop, gather force, and burst into generalized attention. In effect, the determination of the pertinent time-scale is one of the fundamental limitations of the analysis of data from social media, namely to be considered for the posterior processing of information. STE captures the intrinsic time-scale of the information flow and allows a proper diagnosis of the granularity needed to grasp the evolution of social events. Second, we look into the characterization of the structural signature typical of the communication dynamics that underlie social phenomena. We find that the onset of social collective phenomena are characterized by the drop of the characteristic time-scale; we also show that the emergence of coherent patterns of information flow can be mapped into order-disorder transitions in the underlying connectivity patterns of the transfer entropy network. The methodology we present here can therefore be used to gain new insights on the structural and functional relations occurring in large-scale structured populations, eventually leading to the identification of metrics that might be used for the definition of precursors of large-scale social events.

Results

We consider the dataset concerning the time stamped and geolocalized time-series of tweets associated to the following events: the Spanish 15M social unrest in 2011; the *Outono Brasileiro* (“Brazilian Autumn”) in 2013; the discovery of the Higgs boson in 2012; the release of an Hollywood blockbuster in 2012; the acquisition of Motorola by Google in 2011.

[Figure 1 around here]

The spatio-temporal annotation of each tweet in the time series allows the construction of spatially localized activity maps that help identify, as time unfolds, the role that different geographical sub-units played in the global exchange of information. For each dataset the definition of the corresponding spatial unit is performed according to administrative and geographical boundaries as specified in the Material and Methods section (see also Figure 1). Note that the map only shows how the signal increases in all

regions and does not provide evidence of any unexpected transition, pointing out that volume alone is not a good indicator of the evolution of the events.

The time-stamped series of tweets originated from each spatial sub-unit X (supra-urban aggregates) defines the activity time series X_t of the corresponding sub-unit in the social system. Timestamps are modified for each dataset to account for different time zones, see Supplementary Materials (SM) for details. Activity time series encode the role of each geographical sub-unit, a sort of *who-steers-whom*, and several techniques can be used to detect directed exchange of information across the social system. Here, we characterize the dominating direction of information flow between spatial sub-units using Symbolic Transfer Entropy (STE) [21,22]. This well-established technique has been used to infer directional influence between dynamical systems [23-25] and to analyze patterns of brain connectivity [26].

Symbolic transfer entropy quantifies the directional flow of information between two time series X and Y by, first, categorizing the signals in a small set of symbols or alphabet (see Section B.3 of the SM); and, then, computing from the relative frequency of symbols in each sequence \hat{X} and \hat{Y} the joint and conditional probabilities of the sequences indices as

$$T_{Y,X} = \sum p(\hat{x}_{i+\delta}, \hat{x}_i, \hat{y}_i) \log_2 \left(\frac{p(\hat{x}_{i+\delta} | \hat{x}_i, \hat{y}_i)}{p(\hat{x}_{i+\delta} | \hat{x}_i)} \right) \quad (1)$$

where the sum runs over each symbol in the sequence, and $\delta = 1$. The transfer entropy refers to the deviations of the cross Markovian property of the series (independence between them), measured as the Kullback-Leibler divergence [27] (see the SM for all technical details). An important feature of symbolic approaches is that it discounts the relative magnitude of each time series; this is important in our case because different geographical units differ largely in population density or internet penetration rates. Flattening and discretizing the original signal is a key feature to enhance the sensitivity of our proposal to any –even minor– changes in the sub-systems interaction dynamics, see section F in SM.

Within this framework, we first analyze the temporal patterns characterizing the flow of information. Admittedly, micro-blogging data can be sampled at different time-scales Δt . In order to select the optimal sampling rate we consider all possible pairs (X, Y) of geographical units and measure the total STE in the system $T = \sum_{XY} T_{X,Y}$ as a function of Δt . We consider the system-wide characteristic sampling time-scale \mathcal{T} as that which maximizes the total information flow T . This quantity provides an indication of the time-scale at which the information is being exchanged in the system, not necessarily correlated with volume, see early stages of the panels A-D in Figure 2, where, in a very scarce volume scenario, the time scale drops by ~50% (see also Section F in SM). Interestingly, the characteristic time-scale \mathcal{T} changes as the phenomena under analysis unfold, i.e. it decreases as the system approaches the exponential increase in overall activity that signals the onset of the collective phenomena. As shown in the top panels of Figure 2, \mathcal{T} is a proxy for the internally generated coordination in the system that culminates at the very same time of the occurrence of the social event: the street protest day, in the case of political unrest; the movie release date, in the case of the Hollywood blockbuster; and the announcement to the press of the Higgs boson discovery. The only clear exception to this behavior is offered by the company acquisition dataset: the Google-Motorola announcement is a clear example of collective phenomena that is driven mostly by an exogenous factor, i.e. a media announcement. In this case, the dynamical time-scale is constant until the announcement is made public. In the SM we present the same analysis for the randomized signals, showing that time-scale variations are, as expected, washed out from the signal.

The maximized information exchange can be analyzed at the level of geographical subunits by constructing the effective directed network [28] of information flow on a daily basis. This network is

encoded in the matrix $\{T_{XY}\}$ that contains pairwise information about how each component in the system controls (or is controlled by) the others. The matrix $\{T_{XY}\}$ is asymmetric. The directionality is crucial and captures that the geographic area x can exert some driving on area y , and at the same time y might exert some driving on x . For this reason it is convenient to define the directionality index $T_{X,Y}^S = T_{Y,X} - T_{X,Y}$ measuring the balance of information flow in both directions. This index quantifies the dominant direction of information flow and is expected to have positive values for unidirectional couplings with x as the driver and negative values if y is driving x . For symmetric bidirectional couplings we expect $T_{X,Y}^S$ to be null.

[Figure 2 around here]

[Figure 3 around here]

Figure 3 reports the temporal evolution of the maximized $\sum_Y T_{X,Y}^S$ that provides the information flow balance of each specific geographical area. The results show that in the 15M grassroots protests, a limited number of urban areas are initially driving the onset of the social phenomena. These units can mostly be identified with major cities; however the analysis also uncovers *hidden* drivers, such as Orotava, a less known urban area. Only after the first demonstration day on May 15th the driving role becomes much more homogeneously distributed. In the Brazilian case, a set of clear drivers is present only during the onset phase preceding a demonstration on June 6th, becoming fuzzier up to the major demonstration (June 17th) and totally blurred afterwards. We find a similar behavior in the Higgs boson cases (with rumors around the discovery on July 2nd and final announcement on July 4th) [29]. The blockbuster case is driven by a steady excitement of the public before the movie release. Again, as expected, we observe completely different patterns in the case of the Google dataset.

[Figure 4 around here]

In general, the evolving effective networks reveal a transition from a scenario with directed, hierarchical causal relationships to a symmetric *though rather fluctuating* networks where information is flowing symmetrically among all subunits. If information flows mainly in one direction (that is, if the sub-systems are arranged in a highly hierarchical structure) a subunit dominates another, with no or little information flowing in the opposite direction. In this situation, a convenient manipulation of the matrix ($T \rightarrow T^\dagger$) based on a ranking and reordering of the elements according to their directionality index yields an upper triangular matrix (see Materials and Methods). The transition between such hierarchical or centralized driving to a symmetric scenario can be clearly identified monitoring the ratio $\theta = T_l^\dagger / T_u^\dagger$ between the sum all elements of T^\dagger in the lower triangle and the same quantity evaluated in the upper triangle. As schematically illustrated in Figure 4, in a regime of perfect directed driving all the elements below the diagonal are zeros, i.e., $\theta \approx 0$. In the opposite situation (i.e. the perfectly symmetric regime) the values below and above the diagonal are comparable, i.e. $\theta \approx 1$. The quantity θ can thus be considered as a suitable order parameter to characterize this order-disorder transition, thus helping to identify and differentiate communication patterns across the subunits of a system.

Figure 5 shows the behavior of the parameter θ as a function of time in our five datasets. In all the cases we initially observe a highly asymmetric effective network, where a few subunits have a dominant directional coupling to the rest of the system and $\theta \ll 1$. As the systems approach the onset date of the collective event, the quantity $T_l^\dagger / T_u^\dagger$ undergoes a quick transition to $\theta \approx 1$ identifying a regime in which the couplings indicate the existence of collective phenomena where all subunits are mutually affecting each other. We see that in four out of the five datasets the system has a clear order-disorder transition occurring in the proximity of the collective event. Interestingly, in the case of the Brazilian protests the

measure significantly increases before the main event (June 17th). Such behavior probably results from the effects of small precursor protests taking place from June 6th onwards. The same behavior is observed in the Higgs boson dataset, given the existing rumors triggered after July 2nd. Once more, the Google dataset behaves in a completely different way, never showing a clear signature of a collective regime for the couplings network. In the SM we report the same analysis using the randomized signal for both the 15M and the Brazil events, and we observe no order-disorder transition. Similarly, no transition exists for the Twitter unfiltered stream case study (also in SM).

All datasets cover a time-span preceding and following the event, and details on data collection, spatial aggregation (including keyword selection and the geolocalization of messages), and sensitivity analysis of the methodology can be found in the Materials and Methods section and in the SM.

Discussion

The mapping of influence networks using an information theoretic approach offers a new lens to analyze the emergence of collective phenomena. Through this lens, we have revealed the existence of a double transition –in the time scales (slow-to-fast) and directional couplings (hierarchical-to-distributed)– in systems that gather around some sort of collective action. Regarding the first, we bring to light that time series analysis should pay attention to the time scales of the underlying dynamical processes, if it is to provide a reliable account of them –a fact that resonates beyond societal analysis. We also uncover the effective network of information flow between spatially defined sub-units of the social system and study the structural changes of the network connectivity pattern as the system goes through different collective states. In addition, the effective network lends itself to further analysis that can lead to the identification of structural hubs, coordinated communities, influence pathways of geographic or cultural characteristics, and geographical sub-units that may have recurrent roles in the onset of social phenomena. The methodology we present here can therefore be used to gain new insights on the structural and functional relations occurring in large-scale structured populations, eventually leading to the identification of metrics that might be used for the definition of precursors of large-scale social events.

Additionally, the methodology presented here opens interesting paths to advance in the analysis of social phenomena and the identification of generative mechanisms; however, this advance should not be conflated with the possibility of forecasting the emergence of social events –but as a natural complement of techniques along this line, like network change-point detection [30]. The evidence we discuss is agnostic with regard to the predictive potential of online networks and micro-blogging platforms. A real predictive approach cannot be disentangled from an automatic selection of the relevant discussion topics. Our analyses use datasets that were already zooming into the right conversation domain and monitoring specific keywords/hashtags in the Twitter stream. We believe, however, that the general methodological framework we put forward is a first step towards a better understanding of the temporal and spatial signatures of large-scale social events. This advancement might eventually inform the development of tools that can help us anticipate the emergence of macroscopic phenomena. In the meantime, our method offers a valuable resource to analyze how information-driven transitions unfold in socially relevant contexts.

Materials and Methods

Data. The first dataset focuses on the Spanish 15M movement, which emerged in 2011 [31,32]. The data cover a dormant period of low micro-blogging activity that is followed by an explosive phase in which the movement gained the attention of the general public and was widely covered by traditional media sources (see Figure 1). The second dataset contains over 2.5 million geolocalized tweets associated to the *Outono Brasileiro* (“Brazilian Autumn”), a set of political protests that emerged in Brazil in June 2013. Similarly to the Spanish case, the Brazilian data include an initial phase of low activity followed by a

gradual escalation towards the high volumes of general attention that accompanied the street protests. The third dataset tracks communication on the discovery of the Higgs boson before and after it was officially announced to the press in July of 2012; this dataset has been used before to assess how rumors spread through online social networks [29]. The fourth dataset contains messages related to the release of Hollywood blockbuster, announced months prior to its premiere to stir momentum amongst the fan base. Finally, we also consider a dataset tracking communication on the acquisition of Motorola by Google, which came as sudden and unexpected news and immediately triggered a high volume of public attention.

Spanish Twitter activity is spatially coarse-grained according to the list of metropolitan areas defined by the European Spatial Planning Observation Network [33]. This process yields 56 aggregated time series: each of them corresponds to a different geographical area. In addition, there is an extra signal that accounts for any activity not included in those areas, i.e. the system is made up of $N = 57$ components. The data from Brazil are aggregated in 97 basins, which correspond roughly to metropolitan areas [34,35]. The data tracking rumors about the Higgs boson are aggregated at the country level, including only the $N = 61$ most active around this topic. Finally, the Motorola-Google and the blockbuster data are classified in 52 U.S. areas: 50 states, plus Washington D.C and Puerto Rico.

[Figure 5 around here]

Order-Disorder Transition. In real datasets the transition between the different scenarios can be visually inspected with a convenient sorting of the rows and columns of the $T_{x,y}$ matrix. We do so in Figure 5 of the main text, ranking each subunit of the system. The rank for a subunit x is assigned according to the number of times x it is dominant over the rest of the subunits. Once the ranking is settled, any $T_{x,y} < \frac{1}{2}T_{x,y}^{max}$ is set to 0 to improve the visual understanding of the figure. We then obtain a transformed matrix, i.e. $T_{x,y} \rightarrow T_{x,y}^\dagger$. Beyond visualization, the sorted matrix gives room to a monitoring measure $\theta = \frac{\sum_{x>y} T_{x,y}^\dagger}{\sum_{x<y} T_{x,y}^\dagger} = \frac{T_{x,y}^\dagger}{T_{x,y}^\dagger}$ (i.e., the ratio between the sums of all the matrix's elements in the lower and upper triangles) which provides a quantification of the state in which the system is (as explained in the main text). For completion, we have also plotted the same figures without threshold (see Section B.6 of the SM).

References and Notes

- [1] Culotta A (2010) Towards detecting influenza epidemics by analyzing Twitter messages (ACM), pp 115–122.
- [2] Ginsberg J, et al. (2009) Detecting influenza epidemics using search engine query data. *Nature* 457:1012–1014.
- [3] Hickmann KS, et al. (2015) Forecasting the 2013-2014 influenza season using Wikipedia. *PLoS Comput Biol* 11:e1004239.
- [4] Chakraborty P, et al. (2014) Forecasting a moving target: Ensemble models for ill case count predictions. *Proceedings of the 2014 SIAM International Conference on Data Mining*. Proceedings. Society for Industrial and Applied Mathematics pp 262–270.
- [5] Asur S, Huberman BA (2010) Predicting the Future with Social Media, *WI-IAT '10* (IEEE Computer Society, Washington, DC, USA), pp 492–499.
- [6] Bollen J, Mao H, Zeng X (2011) Twitter mood predicts the stock market. *Journal of Computational Science* 2:1–8.

- [7] Curme C, Preis T, Stanley HE, Moat HS (2014) Quantifying the semantics of search behavior before stock market moves. *Proceedings of the National Academy of Sciences* 111:11600–11605.
- [8] Tumasjan A, Sprenger TO, Sandner PG, Welpe IM (2010) Predicting elections with twitter: What 140 characters reveal about political sentiment. *ICWSM* 10:178–185.
- [9] Livne A, Simmons MP, Adar E, Adamic LA (2011) The party is over here: Structure and content in the 2010 election. *ICWSM* 11:17–21.
- [10] Ciulla F, et al. (2012) Beating the news using social media: the case study of American idol. *EPJ Data Science* 1:1–11.
- [11] Botta F, Moat HS, Preis T (2015) Quantifying crowd size with mobile phone and twitter data. *Royal Society Open Science* 2:150162.
- [12] Xu J, Lu TC, Compton R, Allen D (2014) in *Social Computing, Behavioral-Cultural Modeling and Prediction*, Lecture Notes in Computer Science, eds. Kennedy W, Agarwal N, Yang S (Springer International Publishing) Vol. 8393, pp 403–411.
- [13] Ramakrishnan N, et al. (2014) 'Beating the News' with EMBERS: Forecasting Civil Unrest Using Open Source Indicators, *KDD '14 (ACM, New York, NY, USA)*, pp 1799–1808.
- [14] Skoric M, Poor N, Achananuparp P, Lim EP, Jiang J (2012) Tweets and votes: A study of the 2011 singapore general election (*IEEE*), pp 2583–2591.
- [15] Sang ETK, Bos J (2012) Predicting the 2011 Dutch senate election results with twitter (*Association for Computational Linguistics*), pp 53–60.
- [16] Gayo-Avello D (2012) "I wanted to predict elections with twitter and all I got was this lousy paper"—a balanced survey on election prediction using Twitter data. *arXiv preprint arXiv:1204.6441*.
- [17] Tufekci Z (2014) Big questions for social media big data: Representativeness, validity and other methodological pitfalls. *arXiv preprint arXiv:1403.7400*.
- [18] Lazer DM, Kennedy R, King G, Vespignani A (2014) The parable of Google flu: traps in big data analysis. *Science* 343:1203–1205.
- [19] Helbing D (2013) Globally networked risks and how to respond. *Nature* 497:51–59. 20.
- [20] Lehmann J, Goncalves B, Ramasco JJ, Cattuto C (2012) Dynamical Classes of Collective Attention in Twitter (*ACM*), p 251.
- [21] Staniek M, Lehnertz K (2008) Symbolic transfer entropy. *Physical Review Letters* 100:158101.
- [22] Bandt C, Pompe B (2002) Permutation entropy: a natural complexity measure for time series. *Physical Review Letters* 88:174102.
- [23] Schreiber T (2000) Measuring information transfer. *Physical Review Letters* 85:461. 24.
- [24] Hlavackova-Schindler K, Palus M, Vejmelka M, Bhattacharya J (2007) Causality detection based on information-theoretic approaches in time series analysis. *Physics Reports* 441:1-46.
- [25] Ni KY, Lu TC (2014) Information dynamic spectrum characterizes system instability toward critical transitions. *EPJ Data Science* 3:1–25.
- [26] Lizier JT, Heinzle J, Horstmann A, Haynes JD, Prokopenko M (2011) Multivariate information-theoretic measures reveal directed information structure and task relevant changes in FMRI connectivity. *Journal of Computational Neuroscience* 30:85–107.

- [27] Kullback S, Leibler RA (1951) On information and sufficiency. *The Annals of Mathematical Statistics* pp 79–86.
- [28] Sporns O, Chialvo DR, Kaiser M, Hilgetag CC (2004) Organization, development and function of complex brain networks. *Trends in Cognitive Sciences* 8:418–425.
- [29] De Domenico M, Lima A, Mougél P, Musolesi M (2013) The anatomy of a scientific rumor. *Scientific Reports* 3.
- [30] Peel L, Clauset A (2015) Detecting change points in the large-scale structure of evolving networks. *Proc. of the 29th International Conference on Artificial Intelligence (AAAI)*, 2914–2920.
- [31] Borge-Holthoefer J, et al. (2011) Structural and dynamical patterns on online social networks: the Spanish May 15th movement as a case study. *PLoS One* 6:e23883.
- [32] González-Bailón S, Borge-Holthoefer J, Rivero A, Moreno Y (2011) The dynamics of protest recruitment through an online network. *Scientific Reports* 1.
- [33] See <http://www.espon.eu>. Accessed April 16th, 2014.
- [34] Balcan D, et al. (2010) Modeling the spatial spread of infectious diseases: The GLObal Epidemic and Mobility computational model. *Journal of Computational Science* 1:132–145.
- [35] Balcan D, et al. (2009) Seasonal transmission potential and activity peaks of the new influenza A(H1N1): a Monte Carlo likelihood analysis based on human mobility. *BMC Medicine* 7:45.
- [36] American regions. http://en.wikipedia.org/wiki/List_of_regions_of_the_United_States. Accessed August 8th, 2015.
- [37] Brazilian states. http://en.wikipedia.org/wiki/States_of_Brazil. Accessed August 8th, 2015.
- [38] Spain's autonomous communities. http://en.wikipedia.org/wiki/Autonomous_communities_of_Spain. Accessed August 8th, 2015.
- [39] Balcan D, Colizza V, Gonçalves B, Hu H, Ramasco JJ, Vespignani A (2009) Multiscale mobility networks and the large scale spreading of infectious diseases. *Proc. Natl Acad. Sci.*, 106:21484–2189.
- [40] Castells M (2013) *Networks of outrage and hope: Social movements in the internet age*. John Wiley & Sons.
- [41] Columbia University; Center for International Earth Science Information Network (CIESIN) and Centro Internacional de Agricultura Tropical (CIAT). *The Gridded Population of the World Version 3 (GPWv3): Population Grids*. Palisades, NY: Socioeconomic Data and Applications Center (SEDAC), Columbia University. <http://sedac.ciesin.columbia.edu/gpw>.
- [42] Columbia University; International Food Policy Research Institute (IFPRI); The World Bank; Center for International Earth Science Information Network (CIESIN) and Centro Internacional de Agricultura Tropical (CIAT). *Global Rural-Urban Mapping Project (GRUMP), Alpha Version: Population Grids*. Palisades, NY: Socioeconomic Data and Applications Center (SEDAC), Columbia University. <http://sedac.ciesin.columbia.edu/gpw>.
- [43] Conover MD, Davis C, Ferrara E, McKelvey K, Menczer F, Flammini A (2013) The geospatial characteristics of a social movement communication network. *PLoS One*, 8(3):e55957.
- [44] GeoNames. Geonames. <http://www.geonames.org/>, Retr. 2012.
- [45] Gerbaudo P (2012) *Tweets and the streets: social media and contemporary activism*. Pluto Press.

- [46] Granger CWJ (1969) Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: Journal of the Econometric Society*, pages 424–438.
- [47] Jones JJ, Bond RM, Fariss CJ, Settle JE, Kramer ADI, Marlow C, Fowler JH (2013) Yahtzee: An anonymized group level matching procedure. *PLoS One*, 8(2):e55760.
- [48] Kennel MB, Brown R, Abarbanel HDI (1992) Determining embedding dimension for phase-space reconstruction using a geometrical construction. *Physical Review A*, 45(6):3403.
- [49] Schinkel S, Marwan N, Kurths J (2007) Order patterns recurrence plots in the analysis of ERP data. *Cognitive Neurodynamics*, 1(4):317–325.
- [50] Schreiber T. TISEAN software. http://www.mpipks-dresden.mpg.de/~tisean/Tisean_3.0.1.
- [51] Schreiber T (1998) Constrained randomization of time series data. *Physical Review Letters*, 80(10):2105–2108.
- [52] Schreiber T, Schmitz A (1996) Improved surrogate data for nonlinearity tests. *Physical Review Letters*, 77(4):635–638.
- [53] Shannon CE, Weaver W (1948) A mathematical theory of communication, 1948. *Bell Syst. Tech. J.*, 27(379):623.
- [54] Soule, S. A. (1997). The Student Divestment Movement in the United States and Tactical Diffusion: The Shantytown Protest. *Social Forces*, 75(3), 855-882.
- [55] Andrews, K. T., & Biggs, M. (2006). The Dynamics of Protest Diffusion: Movement Organisations, Social Networks, and News Media in the 1960 Sit-Ins. *American Sociological Review*, 71, 752-777.
- [56] Givan, R. K., Roberts, K. M., & Soule, S. A. (Eds.). (2010). *The Diffusion of Social Movements: Actors, Mechanisms, and Political Effects*. Cambridge: Cambridge University Press.
- [57] Wang, D. J., & Soule, S. A. (2012). Social Movement Organizational Collaboration: Networks of Learning and the Diffusion of Protest Tactics, 1960-1995. *American Journal of Sociology*, 117(6).
- [58] Carbunar, B. & Potharaju, R. (2012) You unlocked the Mt. Everest badge on Foursquare! Countering location fraud in geosocial networks. *IEEE 9th International Conference on Mobile Adhoc and Sensor Systems (MASS)*, 182-190.

Acknowledgments. General. We thank D. Allen, R. Compton and T-C Lu at HRL Laboratories LLC for assistance with the Brazilian dataset and useful discussions; we also thank A. Lima for sharing the Higgs boson data. **Funding.** AA acknowledges the support of the European Union MULTIPLEX 317532, the Spanish Ministry of Science and Innovation FIS2012-38266-C02-01, and partial financial support from the ICREA Academia and the James S. McDonnell Foundation. YM acknowledges support from MINECO through Grant FIS2011-25167; Comunidad de Aragón (Spain) through a grant to the group FENOL, and by the EC FET-Proactive Project MULTIPLEX (grant 317532). For the analysis of data outside of the United States of America AV and NP acknowledge the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior National Business Center (DoI/NBC) contract number D12PC00285. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoI/NBE, or the United States Government. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. **Author contributions.** All authors contributed equally to this work. **Data availability.** All data needed to evaluate the conclusions

in the paper are present in the paper, the Supplementary Materials and/or <http://www.jbh.cat/data>
Competing interests. The authors declare no competing interests.

Figures

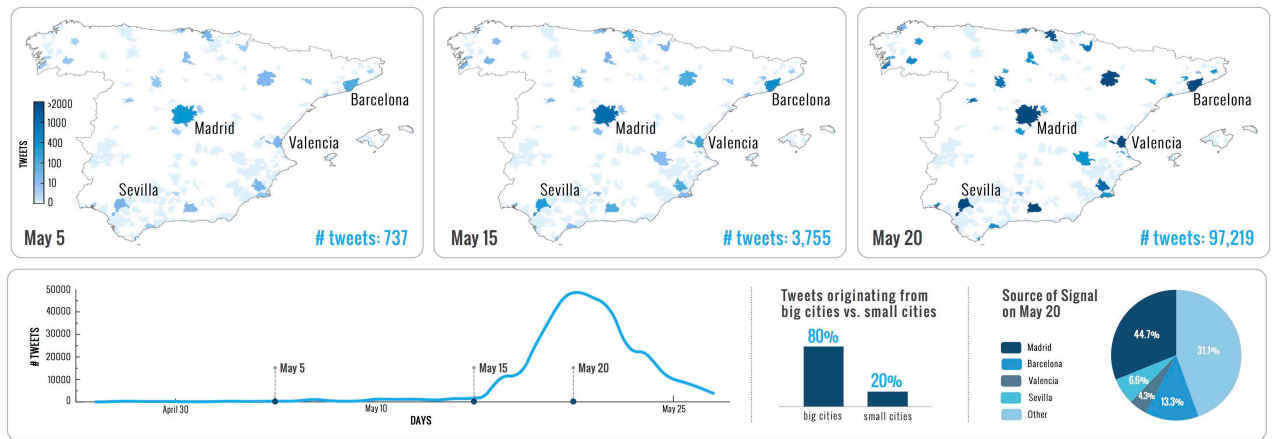


Fig. 1. Spatio-temporal activity as observed from the microblogging platform Twitter. Spain's 15M protest growth in time shows that the protest did not transcend the online sphere until May 15th when the political movement emerged on the streets. Broadcasting traditional media started reporting on it soon after; by that time, demonstrations had been held in the most important cities of the country.

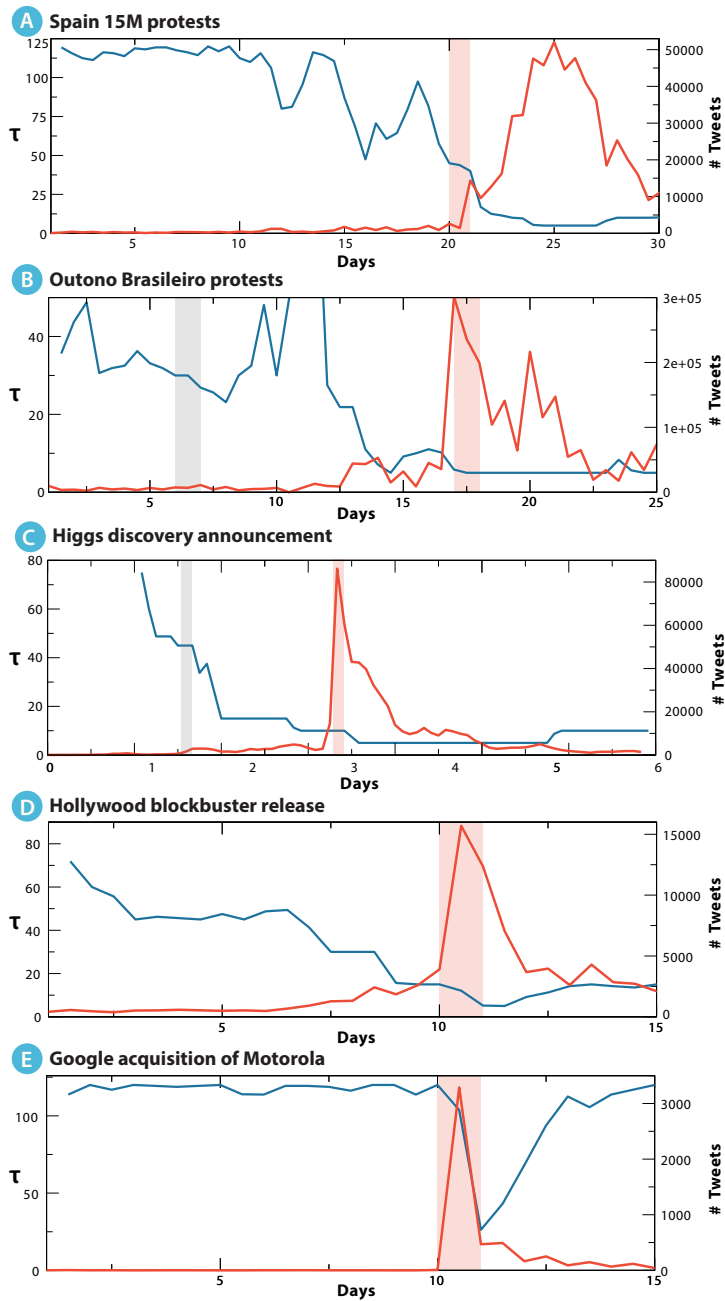


Fig. 2. Characteristic time-scale τ . The panels report the variation of the characteristic time-scale (blue) that maximizes the STE flow as the social event is approached. Red lines correspond to activity volume (number of tweets). Light red vertical lines correspond to the onset of the main social event. Gray ones (in **B** and **C**) indicate a smaller precursor event. (**A**) the 15M event shows a progressive decline of the characteristic time-scale well before the actual social event; the same is observed for the *Outono Brasileiro* in (**B**) (note a data blackout between days 10 and 11). The patterns for the Higgs boson discovery dataset in (**C**) and the Hollywood blockbuster data (**D**) reveal also a drop in the characteristic time-scale, although this is smoother in the movie case. Overall, in all panels (**A-D**) (endogenous activity) the time-scale has dropped already to 50% by the time the absolute volume signals a system-wide event. Finally, (**E**) the Google-Motorola deal triggers a high volume of microblogging activity without actual change in the time-scale of the information flow. In this case the decline is observed in the aftermath of the announcement.

As discussed in the text, this event is the only one that is clearly elicited by an exogenous trigger.

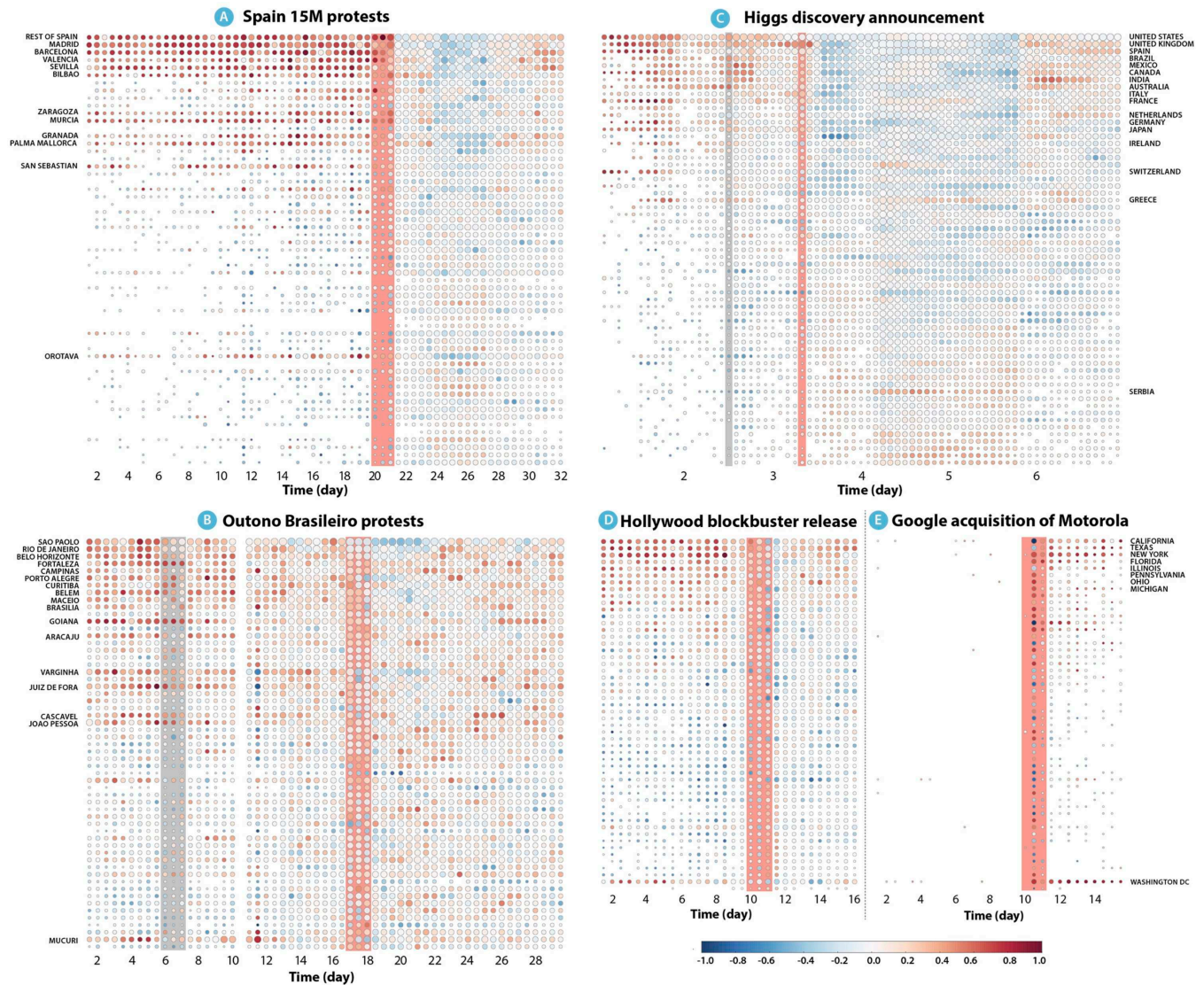


Fig. 3. Evolution of information flow balance between geographical locations for the analyzed events. The color goes from dark blue to dark red (white corresponds to null driving), with the former standing for negative values of $\sum_y T_{X,Y}^S$ (e.g., driven locations) and the latter corresponding to positive information flow balances (i.e., drivers). The size of the circles is log-proportional to the number of messages sent from the location at that time and the vertical bars mark the day of the main event. The geographical locations are ordered according to population size, except for (C), in which countries are ranked with the amount of Higgs-related tweets produced.

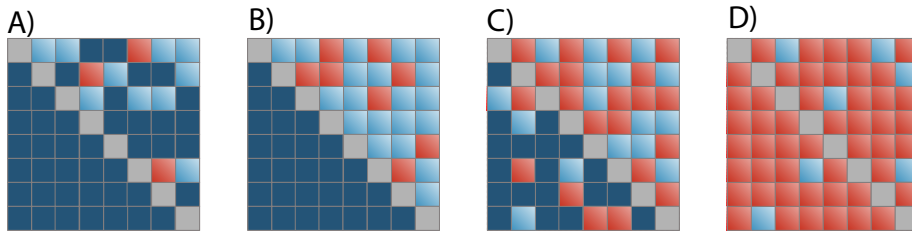


Fig. 4. Schematic representation of a transition from a centralized to a decentralized information flow scenario. If, for any given pair (x, y) , $T_{x,y}^S \sim T_{x,y}$, all existent dynamical driving is net driving, i.e., subsystems present a highly hierarchical structure. In this scenario, if a subsystem dominates another one, the former is not dominated by the latter. This is well illustrated in panels (A) and (B). Note however that in (A), only a few subsystems play an active (dynamical) role; whereas in (B) the situation has reached a perfectly hierarchical structure. Indeed, in this idealized situation the net transfer entropy reaches its maximum: any further addition in terms of dynamical driving will decrease the amount of net transfer entropy (as in panel (C)). Furthermore, (B) and (C) illustrate that there exists a tipping point beyond which the event has necessarily gone global. The extreme case where every subsystem exerts some amount of dynamical driving results in a “null driving” scenario, panel (D). In this schematic representation the color scales goes from dark blue to red, i.e. zero to maximum transfer entropy, respectively.

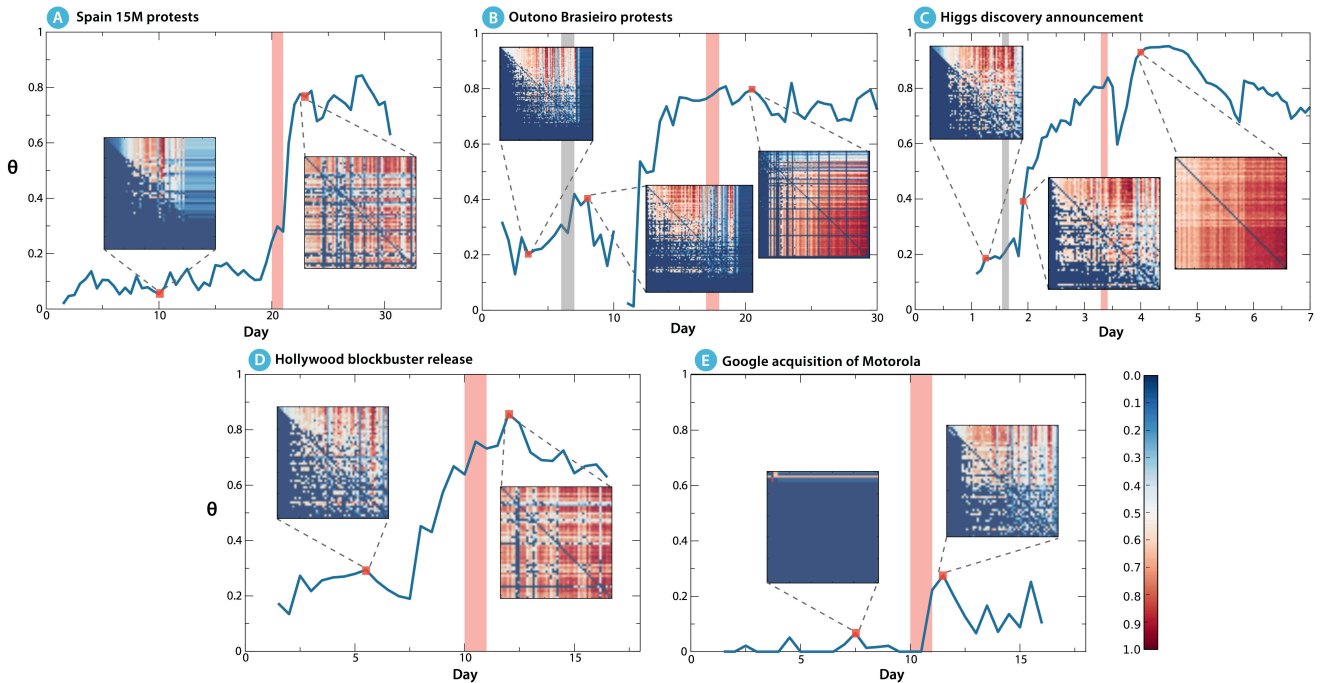


Fig. 5. Order Parameter θ as a function of time for the five events analyzed. The figure represents the behavior of the ratio $\theta = T_l^\dagger / T_u^\dagger$ characterizing the order/disorder of the effective connectivity matrix as a function of time (note a point missing in the Brazilian dataset due to a data blackout between days 10 and 11). For each dataset two or three matrices T^\dagger are plotted considering one or two times before and one after the main event (signaled with a red vertical bar). A clear transition from a hierarchical directed to a distributed symmetrical scenario is observed for the events (A), (B), (C) and (D). The

Google dataset, depicted in panel (E), behaves differently by not showing the same evidence of transition effects.