

Temporal stability in human interaction networks

Renato Fabbri,^{1, a)} Ricardo Fabbri,^{b)} Deborah C. Antunes,^{c)} Marília M. Pisani,^{d)} and Osvaldo N. Oliveira Jr.^{e)}
São Carlos Institute of Physics, University of São Paulo (IFSC/USP), PO Box 369, 13560-970, São Carlos, SP, Brazil

(Dated: 8 March 2017)

This paper reports on stable (or invariant) properties of human interaction networks, with benchmarks derived from public email lists. Activity along time and topology were observed in snapshots in a timeline, and at different scales. Our analysis shows that activity is practically the same for all networks across timescales ranging from seconds to months. The principal components of the participants in the topological metrics space remain practically unchanged as different sets of messages are considered. The activity of participants follows the expected scale-free trace, thus yielding the hub, intermediary and peripheral classes of vertices by comparison against the Erdős-Rényi model. The relative sizes of these three sectors are essentially the same for all email lists and the same along time. Typically, 3-12% of the vertices are hubs, 15-45% are intermediary and 44-81% are peripheral vertices. Similar results for the distribution of participants in the three sectors and for the relative importance of the topological metrics were obtained for 12 additional networks from Facebook, Twitter and Participabr. These properties are consistent with the literature and may be general for human interaction networks, which has important implications for establishing a typology of participants based on quantitative criteria.

PACS numbers: 89.75.Fb, 05.65.+b, 89.65.-s

Keywords: complex networks, pattern recognition, statistics, social network analysis, typology

‘The reason for the persistent plausibility of the typological approach, however, is not a static biological one, but just the opposite: dynamic and social.’ - Adorno et al, 1969, p. 747

I. INTRODUCTION

The first studies dealing explicitly with human interaction networks date from the nineteenth century while the foundation of social network analysis is generally attributed to the psychiatrist Jacob Moreno in mid twentieth century^{1,2}. With the increasing availability of data related to human interactions, research about these networks has grown continuously. Contributions can now be found in a variety of fields, from social sciences and humanities³ to computer science⁴ and physics^{5,6}, given the multidisciplinary nature of the topic. One of the approaches from an exact science perspective is to represent interaction networks as complex networks^{7,8}, with

which several features of human interaction have been revealed. For example, the topology of human interaction networks exhibits a scale-free trace, which points to the existence of a small number of highly connected hubs and a large number of poorly connected nodes. The dynamics of complex networks representing human interaction has also been addressed^{9,10}, but only to a limited extent, since research is normally focused on a particular metric or task, such as accessibility or community detection^{11,12}.

In this paper we analyze the evolution of human interaction networks. Interaction networks from email lists were the most convenient for deriving results and for benchmarking while networks from Facebook, Twitter and Participabr were used for the sake of generalization. Using a timeline of activity snapshots with a constant number of contiguous messages, we found remarkable stability (or invariance) for important network properties. For instance, activity along different timescales follows specific patterns; the most basic topological metrics can always be combined into characteristic principal components; and the fractions of participants in different sectors do not vary with time. This is not an intuitive result, given that participants constantly transition in network structure. Because these properties were shared by networks from various sources, and are consistent with the literature in complex networks¹³, we advocate that the conclusions might be valid for general classes of interaction networks. In particular, this allows us to bridge the gap between data analysis and social sciences in the discussion of types of networks and of participants. It is worth noting that typologies are the canon of scientific literature for the classification of human agents, with pragmatic standards¹⁴ and critical paradigms^{15,16}.

This paper is organized as follows. Section ?? describes related work, while data, scripts and methods of analysis

^{a)} <http://ifsc.usp.br/~fabbri/>; Electronic mail: fabbri@usp.br

^{b)} <http://www.lems.brown.edu/~rfabbri/>; Electronic mail: rfabbri@iprj.uerj.br; Instituto Politécnico, Universidade Estadual do Rio de Janeiro (IPRJ)

^{c)} <http://lattes.cnpq.br/1065956470701739>; Electronic mail: deborahantunes@gmail.com; Curso de Psicologia, Universidade Federal do Ceará (UFC), Campus Sobral

^{d)} <http://lattes.cnpq.br/6738980149860322>; Electronic mail: marilia.m.pisani@gmail.com; Centro de Ciências Naturais e Humanas, Universidade Federal do ABC (CCNH/UFABC)

^{e)} www.polimeros.ifsc.usp.br/professors/professor.php?id=4; Electronic mail: chu@ifsc.usp.br; Also at IFSC-USP

are given in Section ?? and Section ?. Section ? reports results and discussion, leading to Section ? for conclusions. Supplementary data analysis, including directions for video and sound mappings of network structures, and numeric detailed results for networks from Twitter, Facebook and Participabr, are provided in the Supporting Information document.

A. Related work

Research on network evolution is often restricted to network growth, in which there is a monotonic increase in the number of events[?]. Network types have been discussed with regard to the number of participants, intermittence of their activity and network longevity[?]. Two topologically different networks emerged from human interaction networks, depending on whether the frequency of interactions follows a generalized power law or an exponential connectivity distribution[?]. In email list networks, scale-free properties were reported with $\alpha \approx 1.8$ [?] (as in web browsing and library loans[?]), and different linguistic traces were related to weak and strong ties[?].

The fact that unreciprocated edges often exceed 50% in human interaction networks[?] motivated the inclusion of symmetry metrics in our analysis. No correlation of topological characteristics and geographical coordinates was found[?], therefore geographical positions were not considered in our study. Gender related behavior in mobile phone datasets was indeed reported[?] but it is not relevant for the present work because email messages and addresses have no gender related metadata[?].

II. DATA AND SCRIPTS

Email list messages were obtained from the Gmane email archive, which consists of more than 20,000 email lists (discussion groups) and more than 130×10^6 messages[?]. These lists cover a variety of topics, mostly technology-related. The archive can be described as a corpus along with message metadata, including sent time, place, sender name, and sender email address. The usage of the Gmane database in scientific research is reported in studies of isolated lists and of lexical innovations[?] ? .

We observed various email lists and selected four of them together with data from Twitter, Facebook and Participabr for a thorough analysis, from which general properties can be inferred. These lists are as follows:

- Linux Audio Users list[?], with participants from different countries with artistic and technological interests. English is the prevailing language. Abbreviated as LAU from now on.
- Linux Audio Developers list[?], with participants from different countries; a more technical and less active version of LAU. English is the prevailing language. Abbreviated as LAD from now on.

TABLE I. Columns $date_1$ and $date_M$ have dates of first and last messages from the 20,000 messages considered in each email list. N is the number of participants (number of different email addresses), Γ is the number of discussion threads (count of messages without antecedent), \bar{M} is the number of messages missing in the 20,000 collection ($100 \frac{23}{20000} = 0.115$ percent in the worst case).

list	$date_1$	$date_M$	N	Γ	\bar{M}
LAU	2003-06-29	2005-07-23	1147	3374	5
LAD	2003-07-03	2009-10-07	1232	3114	4
MET	2005-08-01	2008-03-07	477	4607	23
CPP	2002-03-12	2009-08-25	1036	4506	7

- Developer's list for the standard C++ library[?], with computer programmers from different countries. English is the prevailing language. Abbreviated as CPP from now on.
- List of the MetaReciclagem project[?], a Brazilian email list for digital culture. Portuguese is the prevailing language, although some messages are written in Spanish and English. Abbreviated as MET from now on.

The first 20,000 messages of each list were considered, with basic attributes of total timespan, authors, threads and missing messages indicated in Table ?. We considered 140 additional email lists to report on the interdependence between the number of participants and the number of discussion threads. Furthermore, 12 networks from Facebook (8), Twitter (2) and Participabr (2) were scrutinized, and their analysis is given in the Supporting Information document for the purpose of testing the generality of the results.

A. Availability

The data and scripts used to derive the results, figures and tables, and this article itself are publicly available. Email messages are downloadable from the Gmane public database[?]. Data annotated from Facebook and Twitter are in a public repository[?]. Data from Participabr were used from the linked data/semantic web RDF triples[?], available in[?]. Computer scripts are delivered through a public domain Python PyPI package and an open Git repository[?]. This open approach to both data and scripts reinforces the scientific aspect of the contribution[?] and mitigates ethical and moral issues involved in researching systems constituted of human individuals[?] ? .

III. METHODS

A. Temporal activity statistics

Messages were counted over time as histograms in the scales of seconds, minutes, hours, days of the week, days of the month, and months of the year. Most standard measures of location and dispersion, e.g. the usual mean and standard deviation, hold little meaning in a compact Riemannian manifold, such as the recurrent time periods that we are interested in. Similar measures were taken using circular statistics[?], in which each measurement t is represented as a unit complex number, $z = e^{i\theta} = \cos(\theta) + i\sin(\theta)$, where $\theta = t\frac{2\pi}{T}$, and T is the period in which the counting is repeated. For example, $\theta = 12\frac{2\pi}{24} = \pi$ for a message sent at $t = 12h$ and given $T = 24h$ for days. The moments m_n , lengths of moments R_n , mean angles θ_μ , and rescaled mean angles θ'_μ are defined as:

$$\begin{aligned} m_n &= \frac{1}{N} \sum_{i=1}^N z_i^n \\ R_n &= |m_n| \\ \theta_\mu &= \text{Arg}(m_1) \\ \theta'_\mu &= \frac{T}{2\pi} \theta_\mu \end{aligned} \quad (1)$$

θ'_μ is used as the measure of location. Dispersion is measured using the circular variance $\text{Var}(z)$, the circular standard deviation $S(z)$, and the circular dispersion $\delta(z)$:

$$\begin{aligned} \text{Var}(z) &= 1 - R_1 \\ S(z) &= \sqrt{-2 \ln(R_1)} \\ \delta(z) &= \frac{1 - R_2}{2R_1^2} \end{aligned} \quad (2)$$

Also, the ratio $r = \frac{b_l}{b_h}$ between the lowest b_l and the highest b_h incidences on the histograms served as a further clue of how close the distribution was to being uniform. As expected, a positive correlation was found in all r , $\text{Var}(z)$, $S(z)$ and $\delta(z)$ dispersion measures, which can be noticed in Section ?? of the Supporting Information. The circular dispersion $\delta(z)$ was found more sensitive and therefore preferred in the discussion of results.

B. Interaction networks

Edges in interaction networks can be modeled both as weighted or unweighted, as directed or undirected^{??}. Networks in this paper are directed and weighted, the most informative of the possibilities. We did not investigate directed unweighted, undirected weighted, and undirected unweighted representations of the interaction networks.

The interaction networks were obtained as follows: a direct response from participant B to a message from participant A yields an edge from A to B, as information went from A to B. The reasoning is: if B wrote a response to a message from A, he/she read what A wrote and formulated a response, so B assimilated information from A, thus $A \rightarrow B$. Edges in both directions are allowed. Each time an interaction occurs, the value of one is added to the edge weight. Selfloops were regarded as non-informative and discarded. Inverting edge direction yields the status network: B read the message and considered what A wrote worth responding, giving status to A, thus $B \rightarrow A$. This paper considers by convention the information network as described above ($A \rightarrow B$) and depicted in Figure ???. These interaction networks are reported in the literature as exhibiting scale-free and small-world^{??} properties, as expected for a number of social networks^{??}.

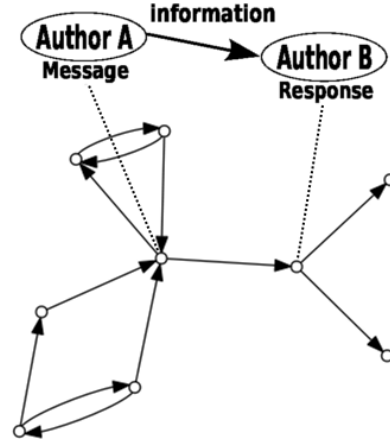


FIG. 1. The formation of interaction networks from exchanged messages. Each vertex represents a participant. A reply message from author B to a message from author A is regarded as evidence that B received information from A and yields a directed edge. Multiple messages add “weight” to a directed edge. Further details are given in Section ??.

1. Topological metrics

The topology of the networks was characterized from a small selection of the most basic and fundamental measurements for each vertex[?], as follows:

- Degree k_i : number of edges linked to vertex i .
- In-degree k_i^{in} : number of edges ending at vertex i .
- Out-degree k_i^{out} : number of edges departing from vertex i .
- Strength s_i : sum of weights of all edges linked to vertex i .

- In-strength s_i^{in} : sum of weights of all edges ending at vertex i .
- Out-strength s_i^{out} : sum of weights of all edges departing from vertex i .
- Clustering coefficient cc_i : fraction of pairs of neighbors of i that are linked, i.e. the standard clustering coefficient metric for undirected graphs.
- Betweenness centrality bt_i : fraction of geodesics that contain vertex i . The betweenness centrality index was computed for weighted digraphs as specified in⁷.

The non-standard metrics below were formulated to capture symmetries in the activity of participants:

- Asymmetry of vertex i : $asy_i = \frac{k_i^{in} - k_i^{out}}{k_i}$.
- Average asymmetry of edges at vertex i : $\mu_i^{asy} = \frac{\sum_{j \in J_i} e_{ji} - e_{ij}}{|J_i|}$, where e_{ij} is 1 if there is an edge from i to j , and 0 otherwise, and J_i is the set of neighbors of vertex i .
- Standard deviation of asymmetry of edges: $\sigma_i^{asy} = \sqrt{\frac{\sum_{j \in J_i} [\mu_i^{asy} - (e_{ji} - e_{ij})]^2}{|J_i|}}$.
- Disequilibrium: $dis_i = \frac{s_i^{in} - s_i^{out}}{s_i}$.
- Average disequilibrium of edges: $\mu_i^{dis} = \frac{\sum_{j \in J_i} \frac{w_{ji} - w_{ij}}{w_{ji} + w_{ij}}}{|J_i|}$, where w_{xy} is the weight of edge $x \rightarrow y$ and zero if there is no such edge.
- Standard deviation of disequilibrium of edges: $\sigma_i^{dis} = \sqrt{\frac{\sum_{j \in J_i} [\mu_i^{dis} - \frac{w_{ji} - w_{ij}}{w_{ji} + w_{ij}}]^2}{|J_i|}}$.

Both standard and non-standard metrics are used for the Erdős sectioning (described in Section ??) and for performing principal component analysis (PCA) (as described in Section ??).

C. Erdős sectioning

It is often useful to think of vertices as hubs, peripheral and intermediary. We have therefore derived the peripheral, intermediary and hub sectors of the empirical networks from a comparison against an Erdős-Rényi network with the same number of edges and vertices, as depicted in Figure ??. We refer to this procedure as *Erdős sectioning*, with the resulting sectors being named as *Erdős sectors*. The Erdős sectioning was recognized as a theoretical possibility by M. O. Jackson in his video lectures⁷, but to our knowledge it has not as yet been applied to empirical data.

The degree distribution $\tilde{P}(k)$ of a real network with a scale-free profile $\mathcal{N}_f(N, z)$ with N vertices and z edges has less average degree nodes than the distribution $P(k)$ of an Erdős-Rényi network with the same number of vertices and edges. Indeed, we define in this work the intermediary sector of a network to be the set of all the nodes whose degree is less abundant in the real network than on the Erdős-Rényi model:

$$\tilde{P}(k) < P(k) \Rightarrow k \text{ is intermediary degree} \quad (3)$$

If $\mathcal{N}_f(N, z)$ is directed and has no self-loops, the probability of the existence of an edge between two arbitrary vertices is $p_e = \frac{z}{N(N-1)}$. A vertex in the ideal Erdős-Rényi digraph with the same number of vertices and edges, and thus the same probability p_e for the presence of an edge, will have degree k with probability

$$P(k) = \binom{2(N-1)}{k} p_e^k (1 - p_e)^{2(N-1)-k} \quad (4)$$

The lower degree fat tail corresponds to the border vertices, i.e. the peripheral sector or periphery where $\tilde{P}(k) > P(k)$ and k is lower than any value of k in the intermediary sector. The higher degree fat tail is the hub sector, i.e. $\tilde{P}(k) > P(k)$ and k is higher than any value of k in the intermediary sector. The reasoning for this classification is as follows: vertices so connected that they are virtually nonexistent in the Erdős-Rényi model, are coherently associated to the hub sector. Vertices with very few connections, which are way more abundant than expected in the Erdős-Rényi model, are assigned to the periphery. Vertices with degree values predicted as the most abundant in the Erdős-Rényi model, near the average, and less frequent in the real network, are classified as intermediary.

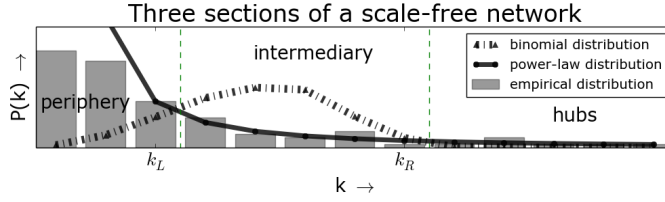


FIG. 2. Classification of vertices by comparing degree distributions[?]. The binomial distribution of the Erdős-Rényi network model exhibits more intermediary vertices, while a scale-free network, associated with the power-law distribution, has more peripheral and hub vertices. The sector borders are defined with respect to the intersections of the distributions. Characteristic degrees are in the compact intervals: $[0, k_L]$, $(k_L, k_R]$, $(k_R, k_{max}]$ for the periphery, intermediary and hub sectors, the “Erdős sectors”. The connectivity distribution of empirical interaction networks, e.g. derived from email lists, can be sectioned by comparison against the associated binomial distribution with the same number of vertices and edges. In this figure, a snapshot of 1000 messages from CPP list yields the degree distribution of an interaction network of 98 nodes and 235 edges. A thorough explanation of the method is provided in Section ??.

To ensure statistical validity of the histograms, bins can be chosen to contain at least η vertices of the real network. The range Δ of incident values of degree k should be partitioned in m parts $\Delta = \cup_{i=1}^m \Delta_i$, with $\Delta_i \cap \Delta_j = \emptyset \forall i \neq j$ and:

$$\Delta_i = \left\{ k \mid \begin{array}{l} \bar{\Delta}_{i-1} < k \leq l \text{ and} \\ \left[\begin{array}{l} N - \sum_{k=0}^{\bar{\Delta}_{i-1}} \eta_k < \eta \text{ and } l = \bar{\Delta} \end{array} \right] \text{ or} \\ \left[\begin{array}{l} \sum_{k=\bar{\Delta}_{i-1}+1}^l \eta_k \geq \eta \text{ and} \\ \left(\sum_{k=\bar{\Delta}_{i-1}+1}^{l-1} \eta_k < \eta \text{ or } l = \bar{\Delta}_{i-1} + 1 \right) \end{array} \right] \end{array} \right\} \quad (5)$$

where η_k is the number of vertices with degree k , while $\bar{\Delta}_{(i)} = \max(\Delta_{(i)})$, and $\bar{\Delta}_0 = -1$. Equation ?? can now be written in the form:

$$\sum_{x=\min(\Delta_i)}^{\bar{\Delta}_i} \tilde{P}(x) < \sum_{x=\min(\Delta_i)}^{\bar{\Delta}_i} P(x) \Leftrightarrow \quad (6)$$

$\Leftrightarrow \Delta_i$ spans intermediary degree values.

If the strength s is used for comparison of the real network against the Erdős-Rényi model, P remains the same, but $P(\kappa_i)$ with $\kappa_i = \frac{s_i}{\bar{w}}$ should be used, where $\bar{w} = 2 \frac{\sum_i s_i}{N}$ is the average weight of an edge and s_i is the strength of vertex i . For in and out degrees (k^{in} , k^{out}),

the real network should be compared against

$$\hat{P}(k^{way}) = \binom{N-1}{k^{way}} p_e^k (1-p_e)^{N-1-k^{way}}, \quad (7)$$

where *way* can be *in* or *out*. In and out strengths (s^{in} , s^{out}) are divided by \bar{w} and compared also using \hat{P} . Note that p_e remains the same, as each edge yields an incoming (or outgoing) edge, and there are at most $N(N-1)$ incoming (or outgoing) edges, thus $p_e = \frac{z}{N(N-1)}$, as with the total degree.

In other words, let γ and ϕ be integers in the intervals $1 \leq \gamma \leq 6$, $1 \leq \phi \leq 3$, and each of the basic six Erdős sectioning possibilities $\{E_\gamma\}$ have three Erdős sectors $E_\gamma = \{e_{\gamma,\phi}\}$ defined as

$$\begin{aligned} e_{\gamma,1} &= \{ i \mid \bar{k}_{\gamma,L} \geq \bar{k}_{\gamma,i} \} \\ e_{\gamma,2} &= \{ i \mid \bar{k}_{\gamma,L} < \bar{k}_{\gamma,i} \leq \bar{k}_{\gamma,R} \} \\ e_{\gamma,3} &= \{ i \mid \bar{k}_{\gamma,i} > \bar{k}_{\gamma,R} \}, \end{aligned} \quad (8)$$

where $\{\bar{k}_{\gamma,i}\}$ is

$$\begin{aligned} \bar{k}_{1,i} &= k_i \\ \bar{k}_{2,i} &= k_i^{in} \\ \bar{k}_{3,i} &= k_i^{out} \\ \bar{k}_{4,i} &= \frac{s_i}{\bar{w}} \\ \bar{k}_{5,i} &= \frac{s_i^{in}}{\bar{w}} \\ \bar{k}_{6,i} &= \frac{s_i^{out}}{\bar{w}} \end{aligned} \quad (9)$$

and both $\bar{k}_{\gamma,L}$ and $\bar{k}_{\gamma,R}$ are found using $P(\bar{k})$ or $\hat{P}(\bar{k})$ as described above and illustrated in Figure ??.

Since different metrics can be used to identify the three types of vertices, more than one metric can be used simultaneously, which is convenient when analysing small networks, such as the cases where only 50 messages are considered in Section ?? of the Supporting Information. After a careful consideration of possible combinations, these were reduced to six:

- **Exclusivist criterion C_1 :** vertices are only classified if the class is the same according to all metrics. In this case, vertices classified do not usually reach N (or 100%), which is indicated by a black line in Figure ??.
- **Inclusivist criterion C_2 :** a vertex has the class given by any of the metrics. Therefore, a vertex may belong to more than one class, and the total number of memberships may exceed N (or 100%), which is indicated by a black line in Figure ??.
- **Exclusivist cascade C_3 :** vertices are only classified as hubs if they are hubs according to all metrics.

Intermediary are the vertices classified either as intermediary or hubs with respect to all metrics. The remaining vertices are regarded as peripheral.

- Inclusive cascade C_4 : vertices are hubs if they are classified as such according to any of the metrics. The remaining vertices are intermediary if they belong to this category for any of the metrics. Peripheral vertices are those which are classified as such with respect to all metrics.
- Exclusivist externals C_5 : vertices are hubs if they are classified as such according to all the metrics. Vertices are peripheral if they are peripheral or hubs for all metrics. The remaining nodes are intermediary.
- Inclusive externals C_6 : hubs are vertices classified as hubs according to any metric. The remaining vertices are peripheral if they are classified as such according to any metric. The rest of the vertices are intermediary.

Using Equations (??), these *compound criteria* C_δ , with δ integer in the interval $1 \leq \delta \leq 6$, can be specified as:

$$\begin{aligned}
C_1 &= \{c_{1,\phi} = \{i \mid i \in e_{\gamma,\phi}, \forall \gamma\}\} \\
C_2 &= \{c_{2,\phi} = \{i \mid \exists \gamma : i \in e_{\gamma,\phi}\}\} \\
C_3 &= \{c_{3,\phi} = \{i \mid i \in e_{\gamma,\phi'}, \forall \gamma, \forall \phi' \geq \phi\}\} \\
C_4 &= \{c_{4,\phi} = \{i \mid i \in e_{\gamma,\phi'}, \forall \gamma, \forall \phi' \leq \phi\}\} \\
C_5 &= \{c_{5,\phi} = \{i \mid i \in e_{\gamma,\phi'}, \forall \gamma, \\
&\quad \forall (\phi' + 1) \% 4 \leq (\phi + 1) \% 4\}\} \\
C_6 &= \{c_{6,\phi} = \{i \mid i \in e_{\gamma,\phi'}, \forall \gamma, \\
&\quad \forall (\phi' + 1) \% 4 \geq (\phi + 1) \% 4\}\}
\end{aligned} \tag{10}$$

Notice that the exclusivist cascade is the same sectioning of an inclusive cascade from periphery to hubs, but with inverted order of sectors. The simplification of all possible compound possibilities to the small set listed above might be formalized in strict mathematical terms, but this was considered out of the scope for current interests.

D. Principal Component Analysis of topological metrics

Principal Component Analysis (PCA) is a well documented technique[?], used here to address the following questions: 1) which metrics contribute to each principal component and in what proportion; 2) how much of the dispersion is concentrated in each component; 3) which are the expected values and dispersions for these quantities over various networks. This enables one to characterize human interaction networks in terms of the relative importance of network metrics and the way they combine.

Let $\mathbf{X} = \{X[i, j]\}$ be a matrix where each element is the value of the metric j at vertex i . Let $\mu_X[j] = \frac{\sum_i X[i, j]}{I}$ be the mean of metric j over all I vertices, $\sigma_X[j] = \sqrt{\frac{\sum_i (X[i, j] - \mu_X[j])^2}{I}}$ the standard deviation of metric j , and $\mathbf{X}' = \{X'[i, j]\} = \left\{ \frac{X[i, j] - \mu_X[j]}{\sigma_X[j]} \right\}$ the matrix with the z -score of each metric. Let $\mathbf{V} = \{V[j, k]\}$ be the matrix $J \times J$ of eigenvectors of the covariance matrix \mathbf{C} of \mathbf{X}' , one eigenvector per column. Each eigenvector combines the original metrics into one principal component, therefore $V'[j, k] = 100 \frac{|V[j, k]|}{\sum_{j'} |V[j', k]|}$ is the percentage of the principal component k that is proportional to the metric j . Let $\mathbf{D} = \{D[k]\}$ be the eigenvalues associated with the eigenvectors \mathbf{V} , then $D'[k] = 100 \frac{D[k]}{\sum_{k'} D[k']}$ is the percentage of total dispersion of the system that the principal component k is responsible for. We consider, in general, the three largest eigenvalues and the respective eigenvectors in percentages: $\{(D'[k], V'[j, k])\}$. These usually sum up between 60 and 95% of the dispersion and reveal patterns for a first analysis. In particular, given L snapshots l of the interaction network, we are interested in the mean $\mu_{V'}[j, k]$ and the standard deviation $\sigma_{V'}[j, k]$ of the contribution of metric j to the principal component k , and the mean $\mu_{D'}[k]$ and the standard deviation $\sigma_{D'}[k]$ of the contribution of the component k to the dispersion of the system:

$$\begin{aligned}
\mu_{V'}[j, k] &= \frac{\sum_{l=1}^L V'[j, k, l]}{L} \\
\sigma_{V'}[j, k] &= \sqrt{\frac{\sum_{l=1}^L (\mu_{V'} - V'[j, k, l])^2}{L}} \\
\mu_{D'}[k] &= \frac{\sum_{l=1}^L D'[k, l]}{L} \\
\sigma_{D'}[k] &= \sqrt{\frac{\sum_{l=1}^L (\mu_{D'} - D'[k, l])^2}{L}}
\end{aligned} \tag{11}$$

The covariance matrix \mathbf{C} is the correlation matrix because \mathbf{X}' is normalized. Therefore, \mathbf{C} is also directly observed as a first clue for patterns by the most simple associations: low absolute values indicate low correlation (and a possible independence); high values indicate positive correlation; negative values with a high absolute value indicate negative correlation. Notice that in this case the variable k is not the degree value but a principal component. In the results the principal components are numbered according to the magnitude of associated eigenvalue and k is incorporated into the notation (e.g. PC2 for metrics of $\mu_{V'}[j, 2]$).

E. Evolution and audiovisualization of the networks

The evolution of the networks was observed within sequences of snapshots. In each sequence, a fixed

TABLE II. The rescaled circular mean θ'_μ and the circular dispersion $\delta(z)$, described in Section ??, for different timescales. This example table was constructed using all LAD messages, and the results are the same for other lists, as shown in Section ?? of the Supporting Information document. The most uniform distribution of activity was found in seconds and minutes. Hours of the day exhibited the most concentrated activity (lowest $\delta(z)$), with mean between 2 p.m. and 3 p.m. ($\theta' = -9.61$). Weekdays, days of the month and months have mean near zero (i.e. near the beginning of the week, month and year) and high dispersion. Note that θ'_u has the dimensional unit of the corresponding time period while $\delta(z)$ is dimensionless.

scale	mean θ'_μ	dispersion $\delta(z)$
seconds	-//-	9070.17
minutes	-//-	205489.40
hours	-9.61	4.36
weekdays	-0.03	29.28
month days	-2.65	2657.77
months	-0.56	44.00

number of messages, i.e. the window size ws , was used for all snapshots. The snapshots were made disjoint in the message timeline, and were used to perform both PCA with topological metrics and Erdős sectioning. Figures and tables were usually inspected with $ws = \{50, 100, 200, 400, 500, 800, 1000, 2000, 2500, 5000, 10000\}$ messages. Variations in the number of vertices, edges and other network characteristics, within the same window size ws , are given in Section ?? of the Supporting Information document.

Network structures were mapped to video animations, sound and musical structures developed for this research[?]. Such *audiovisualizations* were crucial in the initial steps and to guide the research into the most important features of network evolution.

IV. RESULTS AND DISCUSSION

A. Activity along time

Regular patterns of activity were observed along time in the scales of seconds, minutes, hours, days and months. Histograms in each of the time scales were computed as were circular average and dispersion values, and the results are given in Tables ??-??. For example, uniform activity is found with respect to seconds, minutes and days of the months. Weekend days exhibit about half the activity of regular weekdays, and there is a peak of activity between 11am and noon.

In the scales of seconds and minutes, activity is uniform, with the messages being slightly more evenly distributed in all lists than in simulations with the uniform distribution[?]. In the networks, $\frac{\min(\text{incidence})}{\max(\text{incidence})} \in (0.784, .794)$ while simulations reach these values but

TABLE III. Activity percentages along the hours of the day. Nearly identical distributions were observed on other social systems as shown in Section ?? of the Supporting Information document. Highest activity was observed between noon and 6pm (with 1/3 of total day activity), followed by the time period between 6pm and midnight. Around 2/3 of the activity takes place from noon to midnight but the activity peak occurs between 11 a.m. and 12 p.m. This table shows results for the activity in CPP.

	1h	2h	3h	4h	6h	12h
0h	3.66	6.42	8.20	9.30	10.67	33.76
1h	2.76					
2h	1.79	2.88	2.47	3.44	23.09	66.24
3h	1.10					
4h	0.68	1.37	4.35	21.03	28.61	28.61
5h	0.69					
6h	0.83	2.07	18.75	17.59	12.73	8.36
7h	1.24					
8h	2.28	6.80	18.95	25.05	37.63	37.63
9h	4.52					
10h	6.62	14.23	18.68	23.60	28.61	28.61
11h	7.61					
12h	6.44	12.48	12.57	12.58	9.23	4.06
13h	6.04					
14h	6.47	11.02	9.23	8.36	4.30	4.30
15h	6.10					
16h	6.22	12.58	15.88	17.06	17.46	17.46
17h	6.36					
18h	6.01	9.23	12.73	17.59	17.46	17.46
19h	5.02					
20h	4.85	8.36	12.73	17.59	17.46	17.46
21h	4.38					
22h	4.06	8.36	12.73	17.59	17.46	17.46
23h	4.30					

TABLE IV. Activity percentages along weekdays. Higher activity was observed during workweek days, with a decrease of activity on weekend days of at least one third and at most two thirds.

	Mon	Tue	Wed	Thu	Fri	Sat	Sun
LAU	15.71	15.81	15.88	16.43	15.14	10.13	10.91
LAD	14.92	17.75	17.01	15.41	14.21	10.40	10.31
MET	17.53	17.54	16.43	17.06	17.46	7.92	6.06
CPP	17.06	17.43	17.61	17.13	16.30	6.81	7.67

have on average more discrepant higher and lower peaks, i.e. if $\xi = \frac{\min(\text{incidence}')}{\max(\text{incidence}')}$ than $\mu_\xi = 0.7741$ and $\sigma_\xi = 0.02619$. Therefore, the incidence of messages at each second of a minute and at each minute of an hour was considered uniform. In these cases, the circular dispersion is maximized and the mean has little meaning as indicated in Table ??. As for the hours of the day, an abrupt peak is found between 11am and 12pm with the most active period being the afternoon, with one third of total daily activity, and two thirds of activity are allocated in the second 12h of each day. Days of the week revealed a decrease between one third and two thirds of

activity on weekends. Days of the month were regarded as homogeneous with an inconclusive slight tendency of the first week to be more active. Months of the year revealed patterns matching usual work and academic calendars. The time period examined here was not sufficient for the analysis of activity along the years. These patterns are exemplified in Tables ??-??.

B. Stable sizes of Erdős sectors

The distribution of vertices in the hub, intermediary, periphery Erdős sectors is remarkably stable along time if the snapshots hold 200 or more messages, as it is clear in Figure ?? and in Section ?? of the Supporting Information document. Activity is highly concentrated on the hubs, while a very large number of peripheral vertices contribute to only a fraction of the activity. This is expected for a system with a scale-free profile, as confirmed with the distribution of activity among participants in Table ??.

TABLE V. Activity along the days of the month cycle. Nearly identical distributions are found in all systems as indicated in Section ?? of the Supporting Information. Although slightly higher activity rates are found in the beginning of the month, the most important feature seems to be the homogeneity made explicit by the high circular dispersion in Table ??. This specific example and empirical table correspond to the activity of the MET email list.

	1 day	5	10	15 days			
1	3.05	18.25	35.24	50.96			
2	3.38						
3	3.62						
4	4.25						
5	3.94	16.98					
6	3.73						
7	3.17						
8	3.26						
9	3.56	15.73	31.98				
10	3.26						
11	3.81						
12	2.91						
13	3.30	16.25					
14	2.75						
15	2.95						
16	3.36						
17	3.16	15.79	32.78	49.04			
18	3.44						
19	3.36						
20	2.93						
21	3.20	16.99					
22	3.11						
23	3.60						
24	2.74						
25	3.13						
26	3.13						
27	3.07						
28	3.61						
29	3.60						
30	3.57						

TABLE VI. Activity percentages on months along the year. Activity is usually concentrated in Jun-Aug and/or in Dec-Mar, potentially due to academic calendars, vacations and end-of-year holidays. This table corresponds to activity in LAU. Similar results are shown for other lists in Section ?? of the Supporting Information document.

	m.	b.	t.	q.	s.
Jan	10.22	19.56	28.24	35.09	49.16
Fev	9.34				
Mar	8.67	15.53			
Apr	6.86				
Mai	7.28	14.07	20.93	30.36	
Jun	6.80				
Jul	8.97	16.29	24.47		50.84
Ago	7.32				
Set	8.18	16.25		34.55	
Out	8.06				
Nov	7.64	18.30	26.36		
Dez	10.66				

Typically, [3% – 12%] of the vertices are hubs, [15%–45%] are intermediary and [44%–81%] are peripheral, which is consistent with other studies[?]. These results hold for the total, in and out degrees and strengths. Stable sizes are also observed for 100 or less messages if the classification of the three sectors is performed with one of the compound criteria established in Section ??. The networks often hold this basic structure with as few as 10-50 messages, i.e. concentration of activity and the abundance of low-activity participants take place even with very few messages, which is highlighted in Section ?? of the Supporting Information. A minimum window size for the observation of more general properties might be inferred by monitoring both the giant component and the degeneration of the Erdős sectors.

In order to support the generality of these findings, we list the Erdős sector sizes of 12 networks from Facebook, Twitter and Participabr in Table ?? of the Supporting Information document. The fractions of hubs, intermediary and periphery nodes are essentially the same as for the email list networks but with exceptions and a greater variability.

C. Stability of principal components

The principal components of the participants are very stable in the topological space, i.e. in the space of principal components of network measures. Table ?? exemplifies the formation of principal components by providing the averages over non-overlapped activity snapshots of a network. The most important result of this application of PCA, the stability of principal components, is underpinned by the very small dispersion of the contribution of each metric to each principal component.

The first principal component is an average of centrality metrics: degrees, strengths and betweenness central-

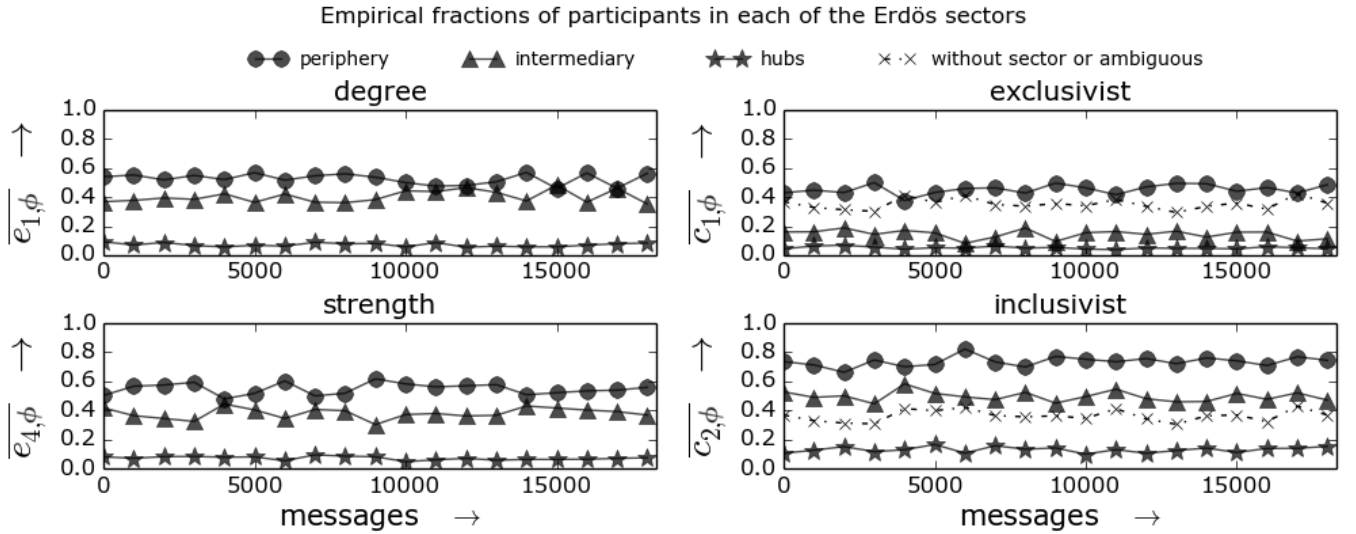


FIG. 3. Stability of Erdős sector sizes. Fractions of participants derived from degree and strength criteria, E_1 and E_4 described in Section ??, are both on the left. Fractions derived from the exclusivist C_1 and the inclusivist C_2 compound criteria are shown in the plots to the right. The ordinates $\overline{e_{\gamma,\phi}} = \frac{|e_{\gamma,\phi}|}{N}$ denote the fraction of participants in sector ϕ through criterion E_γ and, similarly, $\overline{c_{\delta,\phi}} = \frac{|c_{\delta,\phi}|}{N}$ denotes the fraction of participants in sector ϕ through criterion C_δ . Sections ?? and ?? of the Supporting Information bring a systematic collection of such timeline figures with all simple and compound criteria specified in Section ??, with results for networks from Facebook, Twitter and Participabr.

TABLE VII. Distribution of activity among participants. The first column shows the percentage of messages sent by the most active participant. The column for the first quartile (Q_1) gives the minimum percentage of participants responsible for at least 25% of total messages with the actual percentage in parentheses. Similarly, the column for the first three quartiles Q_3 gives the minimum percentage of participants responsible for 75% of total messages. The last decile D_{-1} column shows the maximum percentage of participants responsible for 10% of messages.

list	hub	Q_1	Q_3	D_{-1}
LAU	2.78	1.19 (26.35%)	13.12 (75.17%)	67.32 (-10.02%)
LAD	4.00	1.03 (26.64%)	11.91 (75.18%)	71.14 (-10.03%)
MET	11.14	1.02 (34.07%)	8.54 (75.64%)	80.49 (-10.02%)
CPP	14.41	0.29 (33.24%)	4.18 (75.46%)	83.65 (-10.04%)

ity. On one hand, the similar relevance of all centrality metrics is not surprising since they are highly correlated, e.g. degree and strength have Spearman correlation coefficient $\in [0.95, 1]$ and Pearson coefficient $\in [0.85, 1]$ for window sizes greater than a thousand messages. On the other hand, each of these metrics is related to a different participation characteristic, and their equal relevance for variability, as measured by the principal component, is noticeable. Also, this suggests that these centrality metrics are equally adequate for characterizing the networks and the participants.

According to Table ?? and Figure ??, dispersion is larger in symmetry-related metrics than in clustering co-

TABLE VIII. Loadings for the 14 metrics into the principal components for the MET list, 1000 messages in 20 disjoint positions. The clustering coefficient (cc) appears as the first metric in the table, followed by 7 centrality metrics and 6 symmetry-related metrics. Note that the centrality measurements, including degrees, strength and betweenness centrality, are the most important contributors for the first principal component, while the second component is dominated by symmetry metrics. The clustering coefficient is only relevant for the third principal component. The three components have in average more than 85% of the variance. The low standard deviation σ implies that the principal components are considerably stable.

	PC1		PC2		PC3	
	μ	σ	μ	σ	μ	σ
cc	0.89	0.59	1.93	1.33	21.22	2.97
s	11.71	0.57	2.97	0.82	2.45	0.72
s^{in}	11.68	0.58	2.37	0.91	3.08	0.78
s^{out}	11.49	0.61	3.63	0.79	1.61	0.88
k	11.93	0.54	2.58	0.70	0.52	0.44
k^{in}	11.93	0.52	1.19	0.88	1.41	0.71
k^{out}	11.57	0.61	4.34	0.70	0.98	0.66
bt	11.37	0.55	2.44	0.84	1.37	0.77
asy	3.14	0.98	18.52	1.97	2.46	1.69
μ^{asy}	3.32	0.99	18.23	2.01	2.80	1.82
σ^{asy}	4.91	0.59	2.44	1.47	26.84	3.06
dis	2.94	0.88	18.50	1.92	3.06	1.98
μ^{dis}	2.55	0.89	18.12	1.85	1.57	1.32
σ^{dis}	0.57	0.33	2.74	1.63	30.61	2.66
λ	49.56	1.16	27.14	0.54	13.25	0.95

efficient. We conclude that the symmetry metrics are

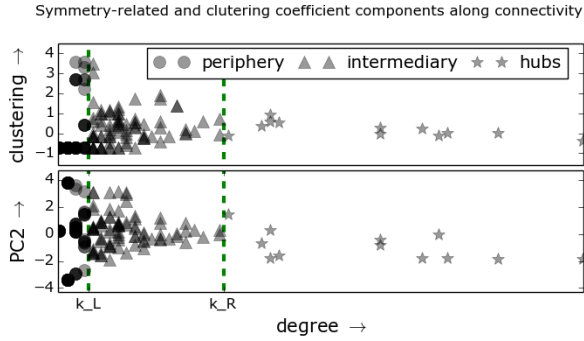


FIG. 4. The first plot highlights the well-known pattern of degree versus clustering coefficient, characterized by the higher clustering coefficient of lower degree vertices. The second plot shows the greater dispersion of the symmetry-related ordinates dominant in the second principal component (PC2). This larger dispersion suggests that symmetry-related metrics are more powerful, for characterizing interaction networks than the clustering coefficient, especially for hubs and intermediary vertices. This figure reflects a snapshot of the LAU list with 1000 contiguous messages.

more powerful, in terms of dispersion in the topological metrics space, in characterizing interaction networks and their participants, than the clustering coefficient, especially for hubs and intermediary vertices (peripheral vertices have larger dispersion with regard to the clustering coefficient). Interestingly, the clustering coefficient is always combined with the standard deviation of the asymmetry and disequilibrium of edges σ^{asy} and σ^{dis} in the third principal component.

Similar results are presented in Sections ?? and ?? of the Supporting Information for other email lists and interaction networks. A larger variability was found for the latter networks, which motivated the use of interaction networks derived from email lists for benchmarking.

D. Types from Erdős sectors

Assigning a type to a participant raises important issues about the scientific cannon for human types and the potential for stigmatization and prejudice. The Erdős sector to which a participant belongs can be regarded as implying a social type for this participant. In this case, the type of a participant changes both along time and as different networks are considered, despite the stability of the network. Therefore, the potential for prejudice of such participant typology is attenuated[?]. In other words, an individual is a hub in a number of networks and peripheral in other networks, and even within the same network he/she most probably changes type along time[?].

The importance of this issue can be grasped by the consideration of static types derived from quantitative criteria. For example, in email lists with a small number of participants, the number of threads has a negative

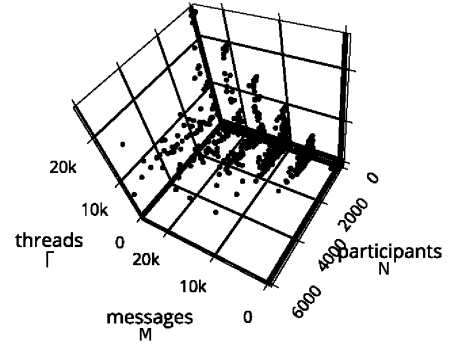


FIG. 5. A scatter plot of number of messages M versus number of participants N versus number of threads Γ for 140 email lists. Highest Γ is associated with low N . The correlation between N and Γ is negative for low values of N but positive otherwise. This negative correlation between N and Γ can also be observed in Table ?? . Accordingly, for $M = 20000$ messages, this inflection of correlation was found around $N = 1500$, while CPP, LAU, LAD, MET lists present smaller networks.

correlation with the number of participants. When the number of participants exceeds a threshold, the number of threads has a positive correlation with the number of participants. This finding is illustrated in Figure ?? and can also be observed in Table ?? . The assignment of types to individuals, in this latter case, has more potential for prejudice because the derived participant type is static and one fails to acknowledge that human individuals are not immutable entities.

Further observations regarding the Erdős sectors and the implicit participant types were made, which are consistent with the literature[?] : 1) hubs and intermediary participants usually have intermittent activity, and stable activity was found only in smaller communities. For instance, the MET list had stable hubs while LAU, LAD and CPP exhibited intermittent hubs. 2) Network structure seems to be most influenced by the activity of intermediary participants as they have less extreme roles than hubs and peripheral participants and can therefore connect to the sectors and other participants in a more selective and explicit manner.

E. Implications of the main findings

The findings reported in this article arose from an exploratory procedure to visually inspect the networks and to analyze considerable amounts of interaction networks data. While this procedure has certainly an ad hoc nature, the statistics in the data are sufficiently robust for important features from these interaction networks to be extracted. Temporal stability, in the sense that interaction networks could be considered as stationary time series, is the most important feature. Also relevant is the significant stability found on the principal compo-

nents, on the fraction of participants in each Erdős Sector and on the activity along different timescales. In fact, these findings confirm our initial hypothesis - based on the literature[?] - that interaction networks should exhibit some stability traces. The potential generality of these findings is suggested by the analysis of networks derived from diverse systems, with interaction networks from public email lists serving as proper benchmarks. Indeed, with such benchmarks one can compare any social network system. Furthermore, this analysis enables us to establish an outline of human interaction networks. It takes the hub, intermediary and periphery sectors out of the scientific folklore and into classes drawn from quantitative criteria. It enables the conception of non-static human types derived from natural properties.

We envisage that the knowledge generated in the analysis may be exploited in applications where the type of each participant and the relative proportion of participants in each sector can be useful metadata. Just by way of illustration, this could be applied in semantic web initiatives, given that the Erdős sectorialization is static in a given snapshot. These results are also useful for classifying resources, e.g. in social media, and for resources recommendation to users[?]. Finally, the knowledge acquired with a quantitative treatment of the whole data may help guide the creation through collective processes of documents to assist in participatory democracy.

Perhaps the most outreaching implications are related to sociological consequences. The results expose a classification of human individuals which is directly related to the concentration of wealth and based on natural laws. The derived human typology changes over different systems and over time in the same system, which implies a negation of the absolute concentration of wealth. Such concentration exists but changes across different wealth criteria and with time. Also, the hubs stand out as dedicated, sometimes enslaved, components of the social system. The peripheral participants have very limited interaction with the network. This suggests that intermediary participants tend to dictate structure, legitimate the hubs and stand out as authorities.

With regard to the limitations of our study, one should emphasize that not all types of human interaction networks were analyzed. Therefore, the plausible generalization of properties has to be treated with caution, as a natural tendency of such systems and not as a rule. Also, the stable properties in the networks were not explored to the limit, which leaves many open questions. For example, what are the maximum and minimum sizes of the networks for which they hold? What is the outcome of PCA analysis when more metrics are considered? What is the granularity in which the activity along the timescales is preserved? Do the findings reported also apply to other systems, beyond human networks?

V. CONCLUSIONS

The very small standard deviations of principal components formation (see Sections ?? and ??), the presence of the Erdős sectors even in networks with few participants (see Sections ?? and ??), and the recurrent activity patterns along different timescales (see Sections ?? and ??), go a step further in characterizing scale-free networks in the context of the interaction of human individuals. Furthermore, the importance of symmetry-related metrics, which surpassed that of clustering coefficient, with respect to dispersion of the system in the topological measures space, might add to the current understanding of key-differences between digraphs and undirected graphs in complex networks. Noteworthy is also the very stable fraction participants in each Erdős sector when the network reaches more than 200 participants. Benchmarks were derived from email list networks and the supplied analysis of networks from Facebook, Twitter and Participabr in the Supporting Information might ease hypothesizing about the generality of these characteristics.

Further work should expand the analysis to include more types of networks and more metrics. The data and software needed to attain these results should also receive dedicated and in-depth documentation as they enable a greater level of transparency and work share, which is adequate for both benchmarking and specifically for the study of systems constituted by human individuals (see Section ??). The derived typology of hub, intermediary and peripheral participants has been applied for semantic web and participatory democracy efforts, and these developments might be enhanced to yield scientific knowledge[?]. Also, we plan to further explore and publish the audiovisualizations used for this research[?] and the linguistic differences found in each of the Erdős sectors[?].

ACKNOWLEDGMENTS

Financial support was obtained from CNPq (140860/2013-4, project 870336/1997-5), United Nations Development Program (contract: 2013/000566; project BRA/12/018) and FAPESP. The authors are grateful to the American Jewish Committee for maintaining an online copy of the Adorno book used on the epigraph[?], to Gmane creators and maintainers for the public email list data, to the communities of the email lists and other groups used in the analysis, and to the Presidency of the Brazilian Republic for keeping Participabr code and data open. We are also grateful to developers and users of Python scientific tools, to Leonardo Paulo Maia (IFSC/USP) and to Francisco J. P. Lopes (UFRJ) for valuable insights.