

Fifteen Minutes of Fame: The Dynamics of Information Access on the Web

Z. Dezső¹, E. Almaas¹, A. Lukács², B. Rácz², I. Szakadát³, A.-L. Barabási¹

1. Center for Complex Network Research and Department of Physics, University of Notre Dame, Notre Dame, IN 46556

2. Computer and Automation Research Institute, Hungarian Academy of Sciences MTA SZTAKI, Budapest, Hungary

3. Axelero Internet Provider Inc., 1364 Budapest, Hungary

(Dated May 13, 2005)

While current studies on complex networks focus on systems that change relatively slowly in time, the structure of the most visited regions of the Web is altered at the timescale from hours to days. Here we investigate the dynamics of visitation of a major news portal, representing the prototype for such a rapidly evolving network. The nodes of the network can be classified into stable nodes, that form the time independent skeleton of the portal, and news documents. The visitation of the two node classes are markedly different, the skeleton acquiring visits at a constant rate, while a news document's visitation peaking after a few hours. We find that the visitation pattern of a news document decays as a power law, in contrast with the exponential prediction provided by simple models of site visitation. This is rooted in the inhomogeneous nature of the browsing pattern characterizing individual users: the time interval between consecutive visits by the same user to the site follows a power law distribution, in contrast with the exponential expected for Poisson processes. We show that the exponent characterizing the individual user's browsing patterns determines the power-law decay in a document's visitation. Finally, our results document the fleeting quality of news and events: while fifteen minutes of fame is still an exaggeration in the online media, we find that access to most news items significantly decays after 36 hours of posting.

The recent interest in the topological properties of complex networks is driven by the realization that understanding the evolutionary processes responsible for network formation is crucial for comprehending the topological maps describing many real systems [1–9]. A much studied example is the WWW, whose topology is driven by its continued expansion through the addition of new documents and links. This growth process has inspired a series of network models that reproduce some of the most studied topological features of the Web [10–17]. The bulk of the current topology driven research focuses on the so called publicly indexable web, which changes only slowly, and therefore can be reproduced with reasonable accuracy. In contrast, the most visited portion of the WWW, ranging from news portals to commercial sites, change within hours through the rapid addition and removal of documents and links. This is driven by the fleeting quality of news: in contrast with the 24-hour news cycle of the printed press, in the online media the non-stop stream of new developments often obliterates an event within hours. But the WWW is not the only rapidly evolving network: the wiring of a cell's regulatory network can also change very rapidly during cell cycle or when there are rapid changes in environmental and stress factors [7]. Similarly, while in social networks the cumulative number of friends and acquaintances an individual has is relatively stable, an individual's contact network, representing those that it interacts with during a given time interval, is often significantly altered from one day to the other. Given the widespread occurrence of these rapidly changing networks, it is important to understand their topology and dynamical features.

Here we take a first step in this direction by studying as a model system a news portal, consisting of news items

that are added and removed at a rapid rate. In particular, we focus on the interplay between the network and the visitation history of the individual documents. In this context, users are often modeled as random walkers, that move along the links of the WWW. Most research on diffusion on complex networks [18–26] ignores the precise *timing* of the visit to a particular web document. There are good reasons for this: such topological quantities as mean free path or probability of return to the starting point can be expressed using the diffusion time, where each time step corresponds to a single diffusion step. Other approaches assume that the diffusion pattern is a Poisson process [27], so that the probability of an HTML request in a dt time interval is pdt . In contrast, here we show that the timing of the browsing process is non-Poisson, which has a significant impact on the visitation history of web documents as well.

I. DATASET AND NETWORK STRUCTURE

Automatically assigned cookies allow us to reconstruct the browsing history of approximately 250,000 unique visitors of the largest Hungarian news and entertainment portal (origo.hu), which provides online news and magazines, community pages, software downloads, free email and search engine, capturing 40% of all internal Web traffic in Hungary. The portal receives 6,500,000 HTML hits on a typical workday. We used the log files of the portal to collect the visitation pattern of each visitor between 11/08/02 and 12/08/02, the number of new news documents released in this time period being 3,908.

From a network perspective most web portals consist of a stable skeleton, representing the overall organization

of the web portal, and a large number of news items that are documents only temporally linked to the skeleton. Each news item represents a particular web document with a unique URL. A typical news item is added to the main page, as well as to the specific news subcategories to which it belongs. For example, the report about an important soccer match could start out simultaneously on the front page, the sports page and the soccer subdirectory of the sports page. As a news document “ages”, new developments compete for space, thus the document is gradually removed from the main page, then from the sports page and eventually from the soccer page as well. After some time (which varies from document to document) an older news document, while still available on the portal, will be disconnected from the skeleton, and can be accessed only through a search engine. To fully understand the dynamics of this network, we need to distinguish between the stable skeleton and the news documents with heavily time dependent visitation.

The documents belonging to the skeleton are characterized by an approximately constant daily visitation pattern, thus the cumulative number of visitors accessing them increases linearly in time. In contrast, the visitation of news documents is the highest right after their release and decreases in time, thus their cumulative visitation reaches a saturation after several days. This is illustrated in Fig. 1, where we show the cumulative visitation for the main page (www.origo.hu/index.html) and a typical news item.

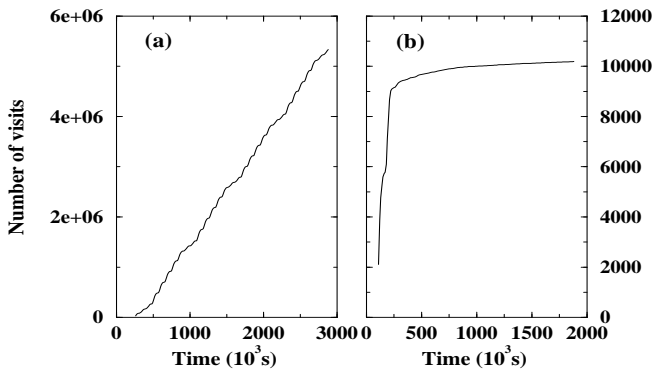


FIG. 1. The cumulative number of visits to a typical skeleton document (a) and a news document (b). The difference between the two visitation patterns allows us to distinguish between news documents and the stable documents belonging to the skeleton.

The difference between the two visitation patterns allows us to distinguish in an automated fashion the web-sites belonging to the skeleton from the news documents. For this we make a linear regression to each site’s cumulative visitation pattern and calculate the deviation from the fitted lines, documents with very small deviations being assigned to the skeleton. The validity of the algorithm was checked by inspecting the URL of randomly selected documents, as the skeleton and the news docu-

ments in most cases have a different format. But given some ambiguities in the naming system, we used the visitation based distinction to finalize the classification of the documents into skeleton and news.

When visiting a news portal, we often get the impression that it has a hierarchical structure. As shown in Fig. 2 the skeleton forms a complex network, driving the visitation patterns of the users. Indeed, the main site, shown in the center, is the most visited, and the documents to which it directly links to also represent highly visited sites. In general (with a few notable exceptions, however), the further we go from the main site on the network, the smaller is the visitation. The skeleton of the studied portal has 933 documents with an average degree close to 2 (i.e. it is largely a tree, with only a few loops, confirming our impression of a hierarchical topology), the network having a few well connected nodes (or hubs), while many are linked to the skeleton by a single link [16,17].

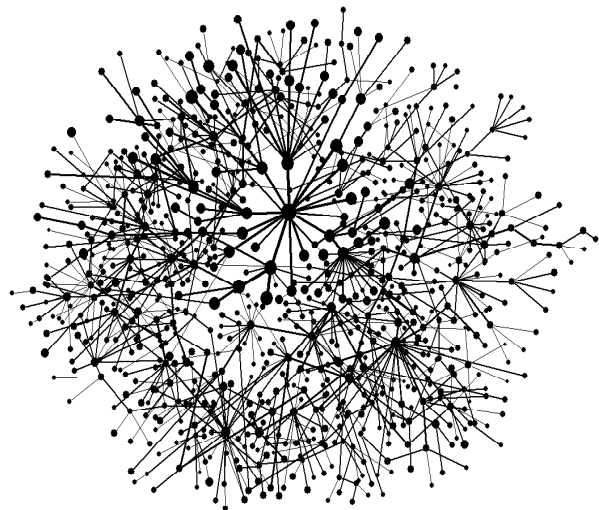


FIG. 2. The skeleton of the studied web portal has 933 nodes. The area of the circles assigned to each node in the figure is proportional with the logarithm of the total number of visits to the corresponding web document. The width of the links are proportional with the logarithm of the total number of times the hyperlink was used by the surfers on the portal. The central largest node corresponds to the main page (www.origo.hu/index.html) directly connected to several other highly visited sites.

II. THE DYNAMICS OF NETWORK VISITATION

Given that the difference between the skeleton and the news documents is driven by the visitation patterns, next we focus on the interplay between the visitation pattern of individual users and the overall visitation of a document. The overall visitation of a specific document is

expected to be determined both by the document’s position on the web page, as well as the content’s potential importance for various user groups. In general the number of visits $n(t)$ to a news document follows a dampened periodic pattern: the majority of visits (28%) take place within the first day, decaying to only 7% on the second day, and reaching a small but apparently constant visitation beyond four days (Fig 3a). Given that after a day or two most news are archived, the long-term saturation of visitation corresponds to direct search or traffic from outside links.

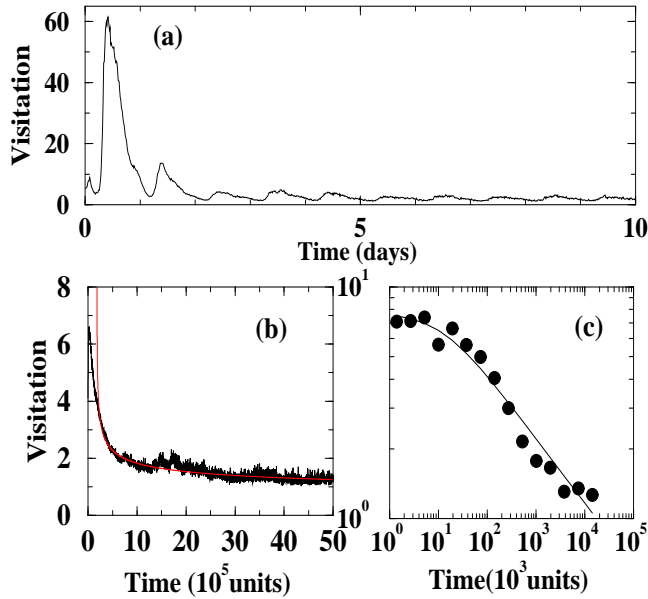


FIG. 3. (a) The visitation pattern of news documents on a web portal. The data represents an average over 3,908 news documents, the release time of each being shifted to day one, keeping the release hour unchanged. The first peak indicates that most visits take place on the release day, rapidly decaying afterward. (b) The same as plot (a), but to reduce the daily fluctuations we define the time unit as one web page request on the portal. (c) Logarithmic binned decay of visitation of (b) shown in a log-log plot, indicating that the visitation follows $n(t) \sim (t + t_0)^{-\beta}$, with $t_0 = 12$ and $\beta = 0.3 \pm 0.1$ shown as a continuous line on both (b) and (c).

To understand the origin of the observed decay in visitation, we assume that the portal has N users, each reading the news document of direct interest for him/her. Therefore, at every time step each user reads a given document with probability p . Users will not read the same news more than once, therefore the number of users which have not read a given document decreases with time. We can calculate the time dependence of the number of potential readers to a news document using

$$\frac{d\mathcal{N}(t)}{dt} = -\mathcal{N}(t)p \quad (1)$$

where $\mathcal{N}(t)$ is the number of visitors which have not read the selected news document by time t . The probability

that a new user reads the news document is given by $\mathcal{N}(t)p$. Equation (1) predicts that

$$\mathcal{N}(t) = N \exp(-t/t_{1/2}) \quad (2)$$

where $t_{1/2} = 1/p$, characterizing the halftime of the news item. The number of visits (n) in unit time is given by

$$n(t) = -\frac{d\mathcal{N}}{dt} = \frac{N}{t_{1/2}} \exp(-t/t_{1/2}). \quad (3)$$

Our measurements indicate, however, that in contrast with this exponential prediction the visitation does not decay exponentially, but its asymptotic behavior is best approximated by a power law (Fig 3c)

$$n(t) \sim t^{-\beta} \quad (4)$$

with $\beta = 0.3 \pm 0.1$, so that while the bulk of the visits takes place at small t , a considerable number of visits are recorded well beyond the document’s release time.

Next we show that the failure of the exponential model is rooted in the uneven browsing patterns of the individual users. Indeed, Eqs. (1) and (2) are valid only if the users visit the site in regular fashion such that they all notice almost instantaneously a newly added news document. In contrast, we find that the time interval between consecutive HTML requests by the same visitor is not uniform, but follows a power law distribution, $P(\tau) \sim \tau^{-\alpha}$, with $\alpha = 1.2 \pm 0.1$ (Fig 4a).

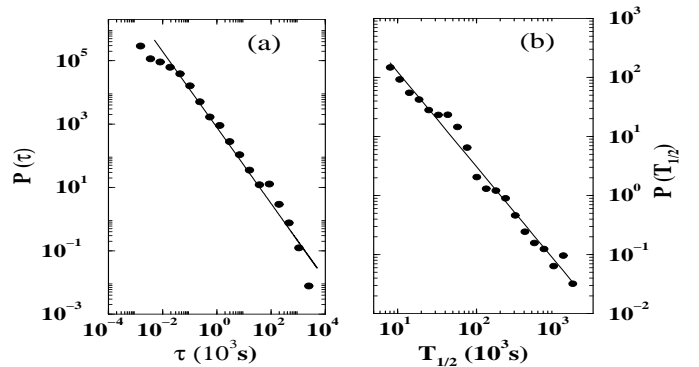


FIG. 4. (a) The distribution of time intervals between two consecutive visits of a given user. The cutoff for high τ ($\tau \approx 10^6$) captures finite size effects, as time delays over a week are undercounted in the month long dataset. The continuous line has slope $\alpha = 1.2$ (b) The halftime distribution for individual news items, following a power-law with exponent -1.5 ± 0.1 .

This means that for each user numerous frequent downloads are followed by long periods of inactivity, a bursting, non-Poisson activity pattern that is a generic feature of human behavior [28] and it is observed in many natural and human driven dynamical processes [29–40]. In the following we show that this uneven user visitation pattern is responsible for the slow decay in the visitation

of a news document and that $n(t)$ can be derived from the browsing pattern of the individual users.

Let us assume that a given news document was released at time t_0 and that all users visiting the main page after the release read that news. Because each user reads each document only once, the visitation of a given document is determined by the number of *new* users visiting the page where the document is featured.

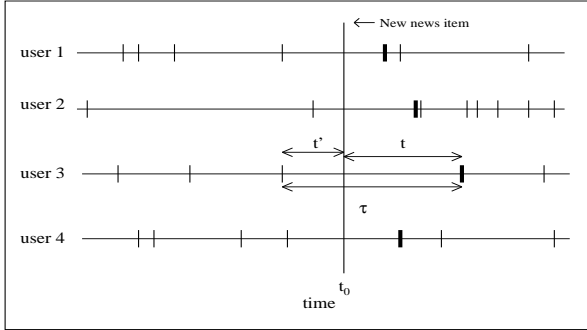


FIG. 5. The browsing pattern of four users, every vertical line representing the time of a visit to the main page. The time a news document was released on the main page is shown at t_0 . The thick vertical bars represent the first time the users visit the main page after the news document was released, i.e. the time they could first visit and read the article.

In Fig. 5 we show the browsing pattern for four different users, each vertical line representing a separate visit to the main page. The thick lines show for each user the first time they visit the main page *after* the studied news document was released at t_0 . The release time of the news (t_0) divides the time interval τ into two consecutive visits of length t' and t , where $t + t' = \tau$. The probability that a user visits at time t after the news was released is proportional to the number of possible τ intervals, which for a given t is proportional to the possible values of t' given by the number of intervals having a length larger than t ,

$$P(\tau > t) = \int_t^\infty \tau^{-\alpha} d\tau \sim t^{-\alpha+1}. \quad (5)$$

If we have N users, each following a similar browsing statistics, the number of new users visiting the main page and reading the news item in a unit time ($n(t)$) follows

$$n(t) \sim NP(\tau > t) \sim Nt^{-\alpha+1}. \quad (6)$$

Equation (6) connects the exponent α characterizing the decay in the news visitation to β in Eq. (4), characterizing the visitation pattern of individual users, providing the relation

$$\beta = \alpha - 1. \quad (7)$$

This is in agreement with our measurements within the error bars, as we find that $\alpha = 1.2 \pm 0.1$ and $\beta = 0.3 \pm 0.1$.

To further test the validity of our predictions we studied the relationship between α and β for the more general case, when a user that visits the main page reads a news item with probability p . We numerically generated browsing patterns for 10,000 users, the distribution for the time intervals between two consecutive visits, $P(\tau)$, following a power-law with exponent $\alpha = 1.5$ (Fig. 6 inset).

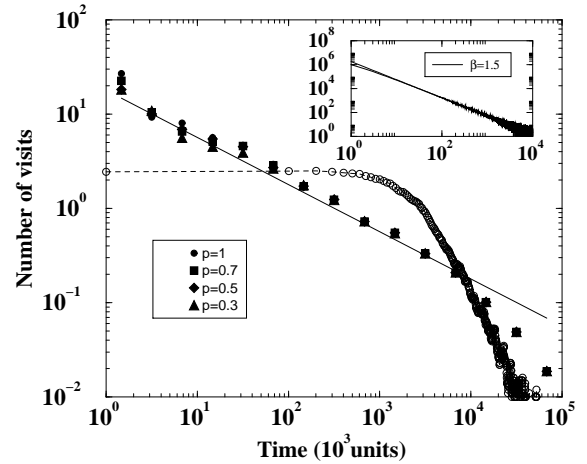


FIG. 6. We numerically generated browsing patterns for 10,000 users, the distribution of the time intervals between two consecutive visits by the same user following a power-law with exponent $\alpha = 1.5$. We assume that users visiting the main page will read a given news item with probability p . The number of visits per unit time decays as a power-law with exponent $\beta = 0.5$ for four different values of p (circles for $p = 1$, squares for $p = 0.7$, diamonds for $p = 0.5$ and triangle for $p = 0.3$). The empty circles represent the visitation of a news item if the users follow a Poisson browsing pattern. We keep the average time between two consecutive visit of each user the same as the one observed in the real data. As the figures indicates, the Poisson browsing pattern cannot reproduce the real visitation decay of a document, predicting a much faster (exponential) decay.

In Fig. 6 we calculate the visits for a given news item, assuming that the users visiting the main page read the news with probability p , characterizing the "stickiness" or the potential interest in a news item. As we see in the figure the value of β is close to 0.5 as predicted by (7). Furthermore, we find that β is independent of p , indicating that the inter-event time distribution $P(\tau)$ characterizing the individual browsing patterns is the main factor that determines the visitation decay of a news document, the difference in the content (stickiness) of the news playing no significant role. As a reference, we also determined the decay in the visitation assuming that the users follow a Poisson visitation pattern [27] with the same inter-event time as observed in the real data. As Fig. 6 shows, a Poisson visitation pattern leads to a much faster decay in document visitation than the power-law seen in Fig. 3c. Indeed, using Poisson inter-event time distribution in (5) would predict an exponentially decaying tail for $n(t)$.

It is useful to characterize the interest in a news document by its half time ($T_{1/2}$), corresponding to the time frame during which half of all visitors that eventually access it have visited. We find that the overall half-time distribution follows a power law (Fig. 4b), indicating that while most news have a very short lifetime, a few continue to be accessed well beyond their initial release. The average halftime of a news document is 36 hours, i.e. after a day and a half the interest in most news fades. A similar broad distribution is observed when we inspect the total number of visits a news document receives (Fig. 7), indicating that the vast majority of news generate little interest, while a few are highly popular [41]. Similar weight distributions are observed in a wide range of complex networks [42–46].

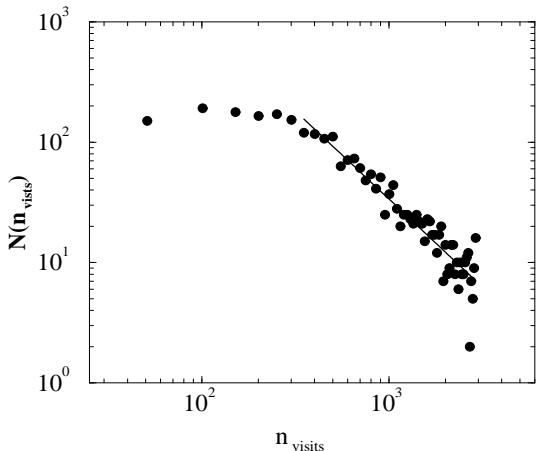


FIG. 7. The distribution of the total number of visits different news documents receive during a month. The tail of the distribution follows a power law with exponent 1.5.

The short display time of a given news document, combined with the uneven visitation pattern indicates that users could miss a significant fraction of the news by not visiting the portal when a document is displayed. We find that a typical user sees only 53% of all news items appearing on the main page of the portal, and downloads (reads) only 7% of them. Such shallow news penetration is likely common in all media, but hard to quantify in the absence of tools to track the reading patterns of individuals.

III. DISCUSSION

Our main goal in this paper was to explore the interplay between individual human visitation patterns and the visitation of specific websites on a web portal. While we often tend to think that the visitation of a given document is driven only by its popularity, our results offer a more complex picture: the dynamics of its accessibility is equally important. Indeed, while “fifteen minutes of

fame” does not yet apply to the online world, our measurements indicate that the visitation of most news items decays significantly after 36 hours of posting. The average lifetime must vary for different media, but the decay laws we identified are likely generic, as they do not depend on content, but are determined mainly by the users’ visitation and browsing patterns [28]. These findings also offer a potential explanation of the observation that the visitation of a website decreases as a power law following a peak of visitation after the site was featured in the media [47]. Indeed, the observed power law decay most likely characterizes the dynamics of the *original* news article, which, due to the uneven visitation patterns of the users, displays a power law visitation decay (see eq. (4)).

These results are likely not limited to news portals. Indeed, we are faced with equally dynamic network when we look at commercial sites, where items are being taken off the website as they are either sold or not carried any longer. It is very likely that the visitation of the individual users to such commercial sites also follows a power law interevent time, potentially leading to a power law decay in an item’s visitation. The results might be applicable to biological systems as well, where the stable network represents the skeleton of the regulatory or the metabolic network, indicating which nodes *could* interact [45,7], while the rapidly changing nodes correspond to the actual molecules that are present in a given moment in the cell. As soon as a molecule is consumed by a reaction or transported out of the cell, it disappears from the system. Before that happens, however, it can take place in multiple interactions. Indeed, there is increasing experimental evidence that network usage in biological systems is highly time dependent [48,49].

While most research on information access focuses on search engines [50], a significant fraction of new information we are exposed to comes from news, whose source is increasingly shifting online from the traditional printed and audiovisual media. News, however, have a fleeting quality: in contrast with the 24-hour news cycle of the printed press, in the online and audiovisual media the non-stop stream of new developments often obliterates a news event within hours. Through archives the Internet offers better long-term search-based access to old events than any other media before. Yet, if we are not exposed to a news item while prominently featured, it is unlikely that we will know what to search for. The accelerating news cycle raises several important questions: How long is a piece of news accessible without targeted search? What is the dynamics of news accessibility? The results presented above show that the online media allows us to address these questions in a quantitative manner, offering surprising insights into the universal aspects of information dynamics. Such quantitative approaches to online media not only offer a better understanding of information access, but could have important commercial applications as well, from better portal design to understanding information diffusion [51–53], flow [54] and marketing in the online world.

-
- [1] R. Albert and A.-L. Barabási, *Rev. Mod. Phys.* **74**, 47 (2002).
- [2] S. N. Dorogovtsev and J. F. F. Mendes, *Evolution of Networks: From Biological Nets to the Internet and WWW* (New York: Oxford University Press, 2003).
- [3] R. Pastor-Satorras and A. Vespignani, *Evolution and Structure of the Internet: A Statistical Physics Approach* (Cambridge University Press, February 2004).
- [4] M. E. J. Newman, *SIAM Review* **45**, 167 (2003).
- [5] M. E. J. Newman, A.-L. Barabási and D. J. Watts, *The Structure and Dynamics of Complex Networks* (Princeton University Press, Princeton, 2005).
- [6] Eli Ben-Naim, H. Frauenfelder, Z. Toroczkai *Complex Networks* (Lecture Notes in Physics) (Springer Verlag: October 16, 2004).
- [7] A.-L. Barabási and Z. N. Oltvai, *Nature Rev. Gen.* **5**, 101-113 (2004).
- [8] S. Bornholdt and H. G. Schuster, *Handbook of Graphs and Networks: From the Genome to the Internet* (Germany: Wiley-VCH 2003).
- [9] S. H. Strogatz, *Nature* **410**, 268 (2001).
- [10] R. Albert, H. Jeong, and A.-L. Barabási, *Nature* **401**, 130 (1999).
- [11] B. A. Huberman, L. A. Adamic, *Nature* **401**, 131 (1999).
- [12] J. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan and A. Tomkins, The web as a graph: measurements, models and methods, *Proc. Intl. Conf. on Combinatorics and Computing*, 1-18, (1999).
- [13] D. M. Pennock, G. W. Flake, S. Lawrence, E. J. Glover, C. L. Giles, *PNAS* **99**, 5207 (2002).
- [14] S. N. Dorogovtsev, J. F. F. Mendes and A. N. Samukhin, *Phys. Rev. Lett.* **85**, 4633 (2000).
- [15] B. Kahng, Y. Park and H. Jeong, *Phys. Rev. E* **66**, 046107 (2002).
- [16] A.-L. Barabási and R. Albert, *Science* **286**, 509 (1999).
- [17] A.-L. Barabási, R. Albert and H. Jeong, *Physica A* **272**, 173 (1999).
- [18] J. D. Noh and H. Rieger, *Phys. Rev. E* **69**, 036111 (2004).
- [19] J. D. Noh and H. Rieger, *Phys. Rev. Lett.* **92**, 118701 (2004).
- [20] S. Jespersen, I. M. Sokolov and A. Blumen, *Phys. Rev. E* **62**, 4405 (2000).
- [21] J. Lahtinen, J. Kertész and K. Kaski, *Phys. Rev. E* **64**, 057105 (2001).
- [22] J. Lahtinen, J. Kertész and K. Kaski, *Physica A* **311**, 571 (2002).
- [23] S. A. Pandit and R. E. Amritkar, *Phys. Rev. E* **63**, 041104 (2001).
- [24] E. Almaas, R. V. Kulkarni and D. Stroud, *Phys. Rev. E* **68**, 056105 (2003).
- [25] R. Monasson, *Eur. Phys. J. B* **12**, 555 (1999).
- [26] B. A. Huberman, P. L. T. Pirolli, J. E. Pitkow and R. M. Lukose, *Science* **280**, 95 (1998).
- [27] J. F. C. Kingman, *Poisson Processes* (Clarendon Press, Oxford, 1993).
- [28] A.-L. Barabási, *Nature* (2005, in press).
- [29] J. F. Omori, *Sci. Imp. Univ. Tokyo* **7**, 111 (1895).
- [30] S. Abe and N. Suzuki, *cond-mat/0410123*.
- [31] A. Vazquez and A.-L. Barabási, preprint.
- [32] C. Dewes, A. Wichmann, A. Feldman, *Proceedings of the 2003 ACM SIGCOMM Conference on Internet Measurement (IMC-03)*, Miami Beach, FL, USA, October 27–29 (ACM Press, New York, 2003).
- [33] S. D. Kleban and S. H. Clearwater, *Hierarchical Dynamics, Interarrival Times and Performance*, *Proceedings of SC'03*, November 15-21, 2003, Phoenix, AZ, USA.
- [34] V. Paxson and S. Floyd, *IEEE/ACM Transactions in Networking* **3**, 226 (1995).
- [35] V. Plerou, P. Gopikrishnan, L. A. N. Amaral, X. Gabaix and H. E. Stanley, *Phys. Rev. E* **62**, R3023 (2000).
- [36] J. Masoliver, M. Montero and G. H. Weiss, *Phys. Rev. E* **67**, 021112 (2003).
- [37] T. Henderson and S. Nhatti, *Modelling user behavior in networked games*, *Proc. ACM Multimedia 2001*, Ottawa, Canada, pp 212-220, 30 September–5 October (2001).
- [38] U. Harder and M. Paczuski, <http://xxx.lanl.gov/abs/cs.PF/0412027>.
- [39] H. R. Anderson, *Fixed Broadband Wireless System Design* (Wiley, New York, 2003).
- [40] P. Ch. Ivanov, B. Podobnik, Y. Lee, and H. E. Stanley *Physica A* **299**, 154-160 (2001).
- [41] F. Menczer, *Proc. Natl. Acad. Sci.* **101**, 5261 (2004).
- [42] K.-I. Goh, B. Kahng, and D. Kim, *Phys. Rev. E* **64**, 051903 (2001); K.-I. Goh, E. Oh, H. Jeong, B. Kahng, and D. Kim, *PNAS* **99**, 12588 (2002).
- [43] A. Barrat, M. Barthelemy, R. Pastor-Satorras and A. Vespignani, *PNAS* **101**, 3747 (2004).
- [44] G. Szabo, M. Alava and J. Kertész *Phys. Rev. E* **66**, 026101 (2002).
- [45] E. Almaas, B. Kovacs, T. Vicsek, Z. N. Oltvai and A.-L. Barabási, *Nature* **427**, 839 (2004).
- [46] S.H. Yook, H. Jeong, A.-L. Barabási and Y. Tu, *Phys. Rev. Lett.* **86**, 5835 (2001).
- [47] A. Johansen and D. Sornette, *Physica A* **276**, 338 (2000).
- [48] J. Jing-Dong et al. *Nature* **430**, 88 (2004).
- [49] N. M. Luscombe et al. *Nature* **431**, 308 (2004).
- [50] S. Lawrence and C. L. Giles, *Nature* **400**, 107 (1999); S. Lawrence and C. L. Giles, *Science* **280**, 98 (1998).
- [51] R. Pastor-Satorras and A. Vespignani, *Phys. Rev. Lett.* **86**, 3200-3203 (2001).
- [52] S. Ciliberti and G. Caldarelli and De los Rios P. Pietronero and L. Zhang, *YC. Phys. Rev. Lett.* **85**, 4848 (2000).
- [53] S. Havlin, and D. Ben-Avraham, *Adv. Phys.* **51**, 187 (2002).
- [54] Z. Toroczkai, and K. E. Bassler, *Nature* **428**, 716 (2004).