# Webgraph connectivity and dynamics: Russian research institutions

**Andrey A. Pechnikov**

Institute of Applied Mathematical Research, Karelian Research Center, Russian Academy of Sciences, Russia. Tel: +7 (8142) 78-11-08. E-mail: pechnikov@krc.karelia.ru

**Anthony M. Nwohiri** (Corresponding author)

Department of Information Management Technology, Federal University of Technology, Owerri, Imo State, Nigeria. Tel: +2348145465855. E-mail address: tonybush98@yahoo.com

## Abstract

This research paper proposes a webgraph dynamics model for research institutes based on a webgraph constructed on a set of instants of time and "return back" through removal of multiple hyperlinks. Analysis of the dynamics model, conducted for the webgraph of the Russian Academy of Sciences (RAS) using real data obtained from websites age determination, shows that emergence of new links for each temporary step in the dynamics model of the webgraph cannot be explained only by the well-known principles of preferential attachment and initial attractiveness of vertices. Of much more importance are the so-called 'administrative actions' – target programs implemented by policy makers both at the regional and national level, which are aimed at developing information resources. Investigations revealed there is an almost perfect match between the model dates marking significant changes in the webgraph dynamics model and the real dates when administrative actions came into force. The webgraph dynamics and connectivity of RAS research institutions were shown to depend on administrative actions of policy makers, implemented at both the regional and national level in the form of targeted programs aimed at developing information resources.

## Keywords

## Introduction

Problems associated with the study of the connectivity of webgraphs of thematic web fragments are today quite conventional. At the forefront are web dynamics problems. Such studies are complex particularly due to lack of source data that needs to be collected in a relatively long period of time. A dynamics model of the webgraph of research institutions of the Russian Academy of Sciences (RAS) is presented. The model is based on a webgraph built for a given (finite) instant of time and "return back" through removal of multiple hyperlinks.

Analysis of this model – conducted by using real data obtained from website dating – shows that the main trend in the development of RAS web space features an exponential increase in the number of the edges (arcs) of the webgraph for every next time interval. This in turn expands the maximal connected component through absorption of smaller components and emergence of new small components.

Investigations conducted showed that emergence of new links at each time step in the webgraph dynamics model cannot be explained only by the well-known principles of preferential attachment and initial attractiveness of vertices. Much more important are the so-called '*administrative actions*' – target programs implemented both at the regional and national levels that are aimed at developing information resources. The study showed a virtually perfect match between the model dates that correspond with significant changes in the webgraph dynamics model and the real dates administrative actions came into force.

## Definitions, tools and research object

The paper deals with the web-based resources of academic institutions in Russia, which, prior to the reform of RAS, were research institutions and members of RAS. Some of these research institutions included the Presidium of RAS, branches of Science, regional branches, research centers and research institutions (RAS, 2015). A Russian Government Decree transferred most of the institutions under the authority of a special administrative Government agency, "Agency of Scientific Institutions" (RGD, 2013). But since information for this study was gathered prior to the reform, we will in this paper refer to the institutions as previously known – "RAS research institutions" (or simply "institutions").

We introduce some definitions and concepts that will be needed hereinafter:

*Definition 1.* Website (or site): a set of HTML-pages and web documents connected by internal hyperlinks, and which have a unity of content identified on the web by its domain name.

*Definition 2.* A unique outbound link is a hyperlink from a set of hyperlinks that have the same address and context, which is located on the highest-level page – the level of the home page of a

site is considered the highest. Henceforth, we will consider only unique outbound links. So the words 'unique' and 'outbound' or 'external' will be mostly omitted.

Generally, a webgraph is a directed graph whose nodes (vertices) represent a set of web pages under study, while the edges (arcs) represent a set of hyperlinks interlinking these pages. We will consider all the pages of one website as a whole. Therefore, by aggregating the hyperlinks appropriately, we get a webgraph of a set of websites. Using only unique hyperlinks (in the sense of Definition 2), we get a webgraph built based on a set of websites with unique hyperlinks. This graph has no loops, but there can be multiple (parallel) edges since existence of several hyperlinks between two websites is possible.

The webgraph under study is denoted as $G=G(S,V)$, where $S$ is the set of nodes (representing the official websites of the institutions), and $V$ is a set of edges (representing the hyperlinks connecting the websites). We are investigating 397 official websites of institutions, therefore $|S|=397$.

Data from a project on webometric ranking for research institutions were used as basis for formation of $S$ (WRR, 2015). A complete list of the websites can be found here: http://webometrics-net.ru/doc/science_sites2014.xls. The table below shows only a few of them.

**Table 1. Official websites of RAS institutions**

| S/N | Name | Domain name |
|---|---|---|
| 1 | Russian Academy of Sciences | http://www.ras.ru |
| 2 | Branch of Mathematics, RAS | http://omn.ras.ru |
| 13 | Far Eastern Branch, RAS | http://www.febras.ru |
| 22 | Karelian Research Centre, RAS | http://www.krc.karelia.ru |
| 56 | Institute of Nanotechnology of Microelectronics, RAS | http://www.intm-ras.ru |
| 70 | Institute of Informatics Problems, RAS | http://www.ipiran.ru |
| 297 | Institute of Oil and Gas Problems, Siberian Branch, RAS | http://ipng.ysn.ru |

Specialized software BeeCrawler (Pechnikov & Chernobrovkin, 2014) was used to search for and collect hyperlinks. Database of External Hyperlinks (DBEH) (Golovin *et al.*, 2013) was used to store and process hyperlink data. DBEH can secure guest access on the web. To gain access, the user needs to follow the link "http://grid.krc.karelia.ru/webometrics2" and type "guest" as the username and password. Open software platform Gephi (Gephi, 2015) was employed for webgraph analysis.

In the end, about 18,000 unique hyperlinks interlinking the 397 websites were found. A set of edges $V$ was formed for webgraph $G=G(S,V)$, $|V| = 17,794$. The multiplicity of edges ranges from 1 to 3100. When each multiple edge is considered as one, we have 2331 edges all together. There more than 10 multiplicity of edges only in 63 cases, that is, less than 3% of all cases.

## Connectivity properties of webgraphs of official websites of RAS institutions

We proceed from webgraph $G(S,V)$ to webgraph $G(S,V(n))$, built on the same set of nodes. Set of edges $V(n)$ is obtained from the set $V$ by the following rule:

1.  Let nodes $i,j \in S$ and let there exist a subset of multiple edges $\{(i,j), (i,j), …\} \subseteq V$, connecting $i$ and $j$;

2.  Let parameter $n$ be an integer greater than 0,

3.  Then the edge $(i,j) \in V(n)$, if and only if $|\{(i,j), (i,j), …\}| \geq n$.

In other words, edge $(i,j)$ exists in $G(S,V(n))$, if there are $n$ edges $(i,j)$ in $G(S,V)$ connecting those same nodes. It follows that $G(S,V(n))$ is a directed graph without multiple edges and loops.

Henceforth, we will be considering $n \in 1..10$, since for $n>10$, the corresponding webgraph $G(S, V(n))$ has weak connectivity and its nodes are mostly isolated.

A sequence of web counts was constructed for the previously obtained graph $G(S,V)$:

$$G(S, V(1)), G(S, V(2)), …, G(S, V(10)) \hspace{2cm} (*).$$

Some of these webgraphs are shown in Figure 1 for given values of $n$. Isolated nodes are not shown. This is to make the webgraph easier for the viewer to comprehend. As one would expect, the webgraph connectivity weakens as $n$ increases.
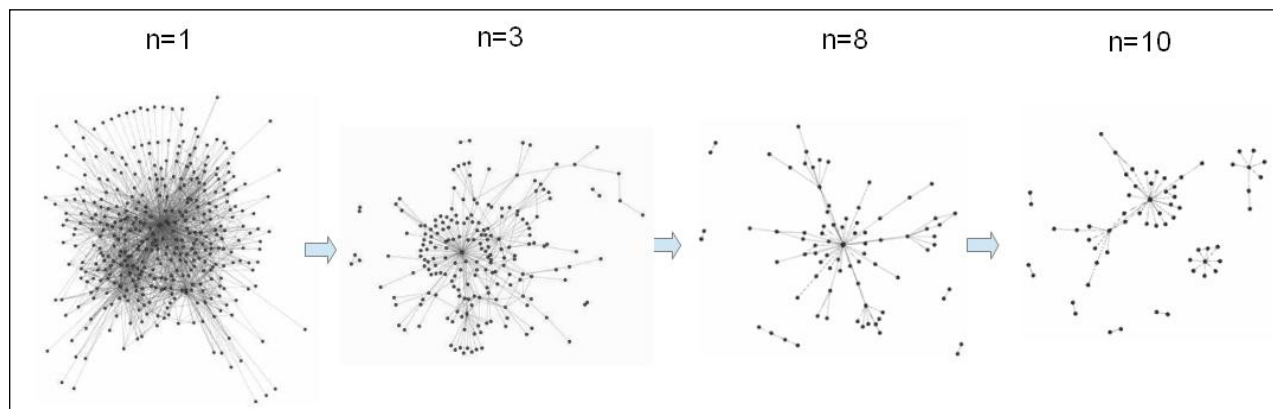


**Figure 1. Changes in webgraph G(S,V(n)) as parameter *n* increases**

The following three webgraph connectivity characteristics were considered: number of nodes having at least one hyperlink incident to them (henceforth referred to as "non-isolated nodes"), total number of edges, and number of nodes in the maximal strongly connected component (SCC) in cases where there are more than one SCC. It is clear that lower number of *n*, which entails lowering the number of hyperlinks linking the websites, leads to an increase in the values of these characteristics. For a sequence of webgraphs (*) constructed based on a set of official websites *S*, changes in the characteristics, depending on *n*, are shown in the form of graphs in Figure 2. The values of *n* on the x-axis are taken in descending order, and this order is used henceforth in this paper.
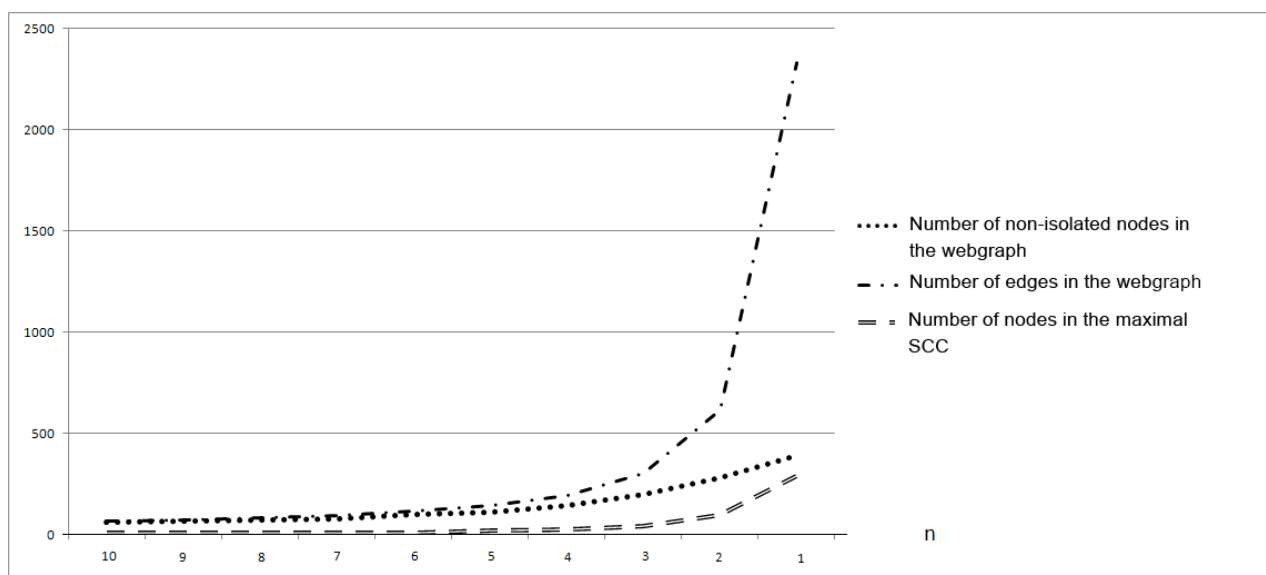


**Figure 2. Connectivity characteristics values for the sequence G(S, V(1)), G(S, V(2)), …, G(S, V(10)) (*)**

The graphs show that the softer the requirements for multiplicity of hyperlinks, the greater the values of all the characteristics, but the order of their growth is different. The number of edges in the webgraph shows exponential growth. The maximum possible number of edges is equal to $|S|*(|S|-1)$, which, in our case, is over 157,000 edges. It means we are still very far away from saturating the graph with edges.

The number of non-isolated nodes grows as a degree 2 polynomial, while the number of nodes in the maximal SCC increases as a degree 3 polynomial. It is obvious that the number of non-isolated nodes and nodes in the maximal SCC in our case are bounded above by $|S|=397$, i.e., for $n=1$, these characteristics are close to the limit.

## Dynamics model of the webgraph of RAS research institutions

We consider the following dynamics model of the webgraph of RAS research institutions, built on the basis of sequence (*), taken in reverse order.

Let $t_1$, $t_2$,...,$t_{10}$ be some conditional discrete instants in time; $t_i < t_j$ for $i < j$. We will assume that webgraph $G(S, V(10))$ reflects the state of the web space of the research institutions at initial time $t_1$, webgraph $G(S, V(9))$ at time $t_2$, etc. From Figures 1 and 2, it is evident that the number of edges in the graph at time $t_i + 1$ is greater than at time $t_i$, and this increase is quite naturally explained by the fact that new hyperlinks interlinking the websites appeared during time interval $(t_i, t_i + 1]$. Such interpretation allows to see how the webgraph changes with increasing number of edges. If one mentally "scrolls" Figure 2 from left to right, one would see how the strongly connected components, the non-isolated and isolated nodes are behaving as time changes.

We examine Figure 3 in more detail. The figure shows transition of the graph from the state at time point $t_1$ to $t_2$, and then to $t_5$ (the states at time $t_3$ and $t_4$ are omitted because significant events did not occur at these instants of time).
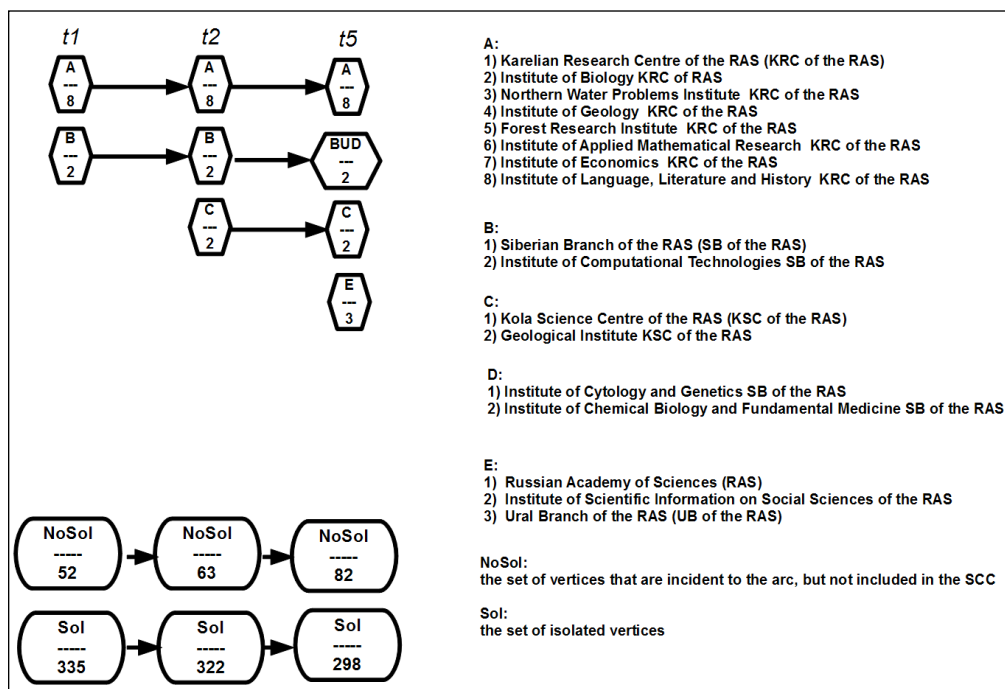


**Figure 3. States of the webgraph at instants of time t1, t2 and t5**

The "virtual" column relating to time $t_1$ shows four disjoint subsets of nodes: *A* (maximal SCC), *B* (one more SCC), *NoSol* (non-isolated nodes) and *Sol* (isolated nodes); $A \cup B \cup NoSol \cup Sol = S$.

The number written under each set is equal to the power of that set. 397 nodes in *G(S,V(10))* have 69 edges only.

The influence of regional dimension was visible in the early stages of the development of the academic web fragment. For example, at instant of time $t_1$, the set *A* consists of the websites of the Karelian Research Center of the Russian Academy of Sciences, and 7 of its daughter institutions; let's call it the "Karelian" component (Republic of Karelia is a federal subject of the Russia). The set *C* is a "Murmansk" component (Murmansk Oblast is the subject of Russia, mainly located on the Kola Peninsula). Set *B* (at instants of time $t_1$ and $t_2$) and set $B \cup D$ (at $t_5$) represent the "Siberian" component.

Let's also note that set *E* appears at time $t_5$, consisting of three sites, including the official website of RAS, which has no regional dimension. It is shown later that the presence of set *E* had a decisive influence on the emergence of a single large component.

Upon transition from $t_5$ to $t_6$ (see Figure 4), five "small" SCCs *A-E* (with addition of a new SCC *F*) are joined to the maximal SCC containing 21 sites. At time $t_6$, the webgraph contains 113 non-isolated nodes and 146 edges.
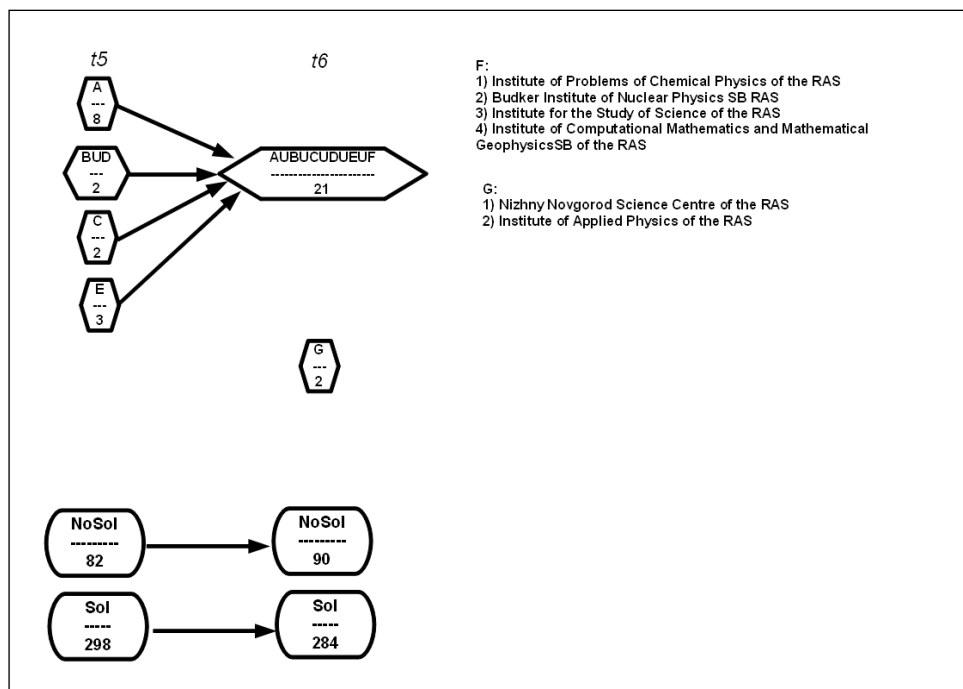


**Figure 4. Changes in the webgraph as time changes from t5 to t6**

Next, Figure 5 shows the states at $t_6$ and $t_8$ (the state at $t_7$ is not fundamentally different from that of $t_6$).
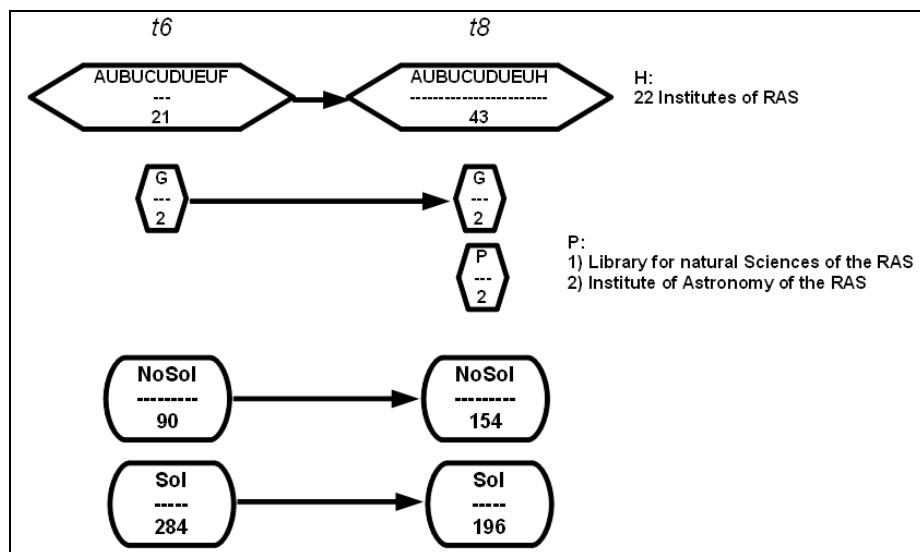
**Figure 5. State of the webgraph at times t6 and t8**

Rapid expansion in the maximum SCC (almost twice) and emergence of new small SCCs have been observed; moreover, the webgraph at time $t_8$ contains 201 non-isolated nodes and 306 edges. The final Figure 6 shows the states of the webgraph at $t_8$, $t_9$ and $t_{10}$. Attention is drawn to the emergence of SCC $R$ at $t_9$, containing 6 websites of the institutions of the Far Eastern Branch of RAS ("Far Eastern" component). Although there are still 4 isolated sites and 106 non-isolated sites outside the SCC, it can be assumed that formation of the main component of the webspace containing over 70% of the sites has been completed.
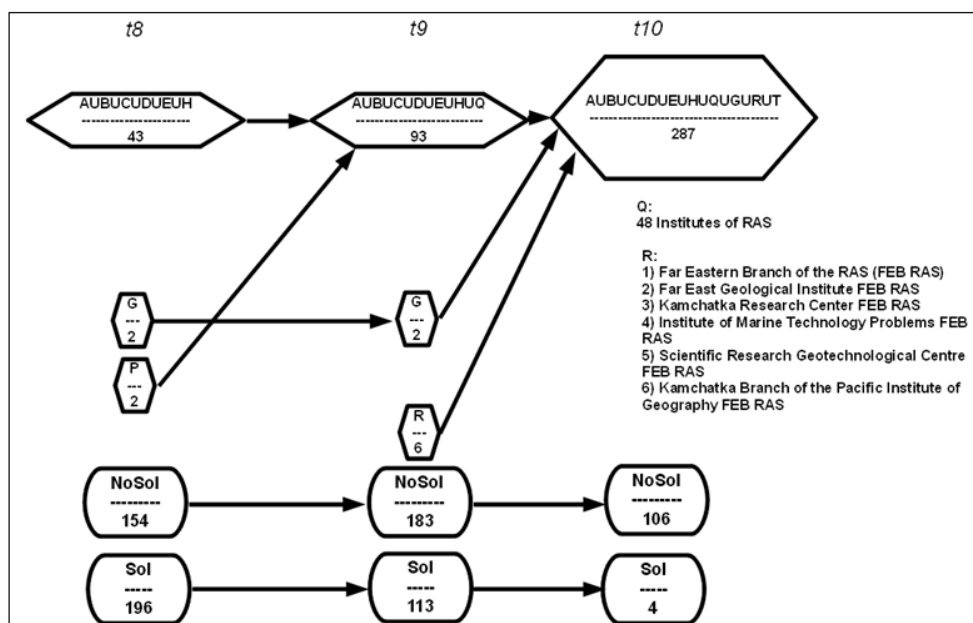


**Figure 6. Final states of the webgraph**

## Results and discussion

The proposed webgraph dynamics model reveals a major trend in the development of the web space of official sites of RAS research institutions: for each step in time, there is an exponential increase in the number of edges (arcs) of the webgraph, which leads to an increase in the maximum component as a result of absorption of smaller components and subsequent emergence of new small components. At the starting point, the model contains two components – Karelian and Siberian components; the Kola component is then added to them. Then, the Siberian component increases slightly, and then emerges a new component that includes the official website of RAS. After this, the main component starts growing rapidly, thus reducing the number of isolated nodes to only 4.

It would seem that this process confirms such webgraph modeling rules as preferential attachment and initial attractiveness of nodes (Bollobas & Riordan, 2003). Let's recall that based on the preferential attachment model, it is believed that a new site appears at each instant of time. The new site contains a fixed number of links pointing to its predecessor sites. Here, the likelihood that the new website will place a link pointing to one of the previous sites is proportional to the number of already existing links to that previous site. Initial attractiveness may be understood as a certain preference for the new site to place links to already existing sites. For example, there will soon be a link pointing to an institution working in the same scientific field as the institution that owns the new site or that is a part of the same research center rather than a link to the website of an institution from other scientific fields.

Using data from the "*Internet Archive: Wayback Machine*" project (Internet Archive, 2015), we attempted to establish the date when the sites of RAS research institutions were created. For example, the first version of the site of the Russian Academy of Sciences (www.ras.ru) dates back to July 1997, Karelian Research Center (www.krc.karelia.ru) to April 1997, and the Institute of Applied Mathematical Research (mathem.krc.karelia.ru) to March 2007. This article does not date all the sites of RAS research institutions. However, conclusions drawn from analysis of this information suggests that emergence of new links at each instant of time in the webgraph dynamics model is impossible to explain only by the rules of preferential attachment and initial attractiveness of nodes. Such a discovered fact that a component containing two websites of institutions working in the same field of science (biology) was observed only once in the dynamics model of the webgraph (see Figure 3, set D) serves as a reason to abandon the principles of preferential attachment and initial attractiveness.

The influence of regional dimension noted earlier is not a confirmation of the principle of initial attractiveness on the basis of geographic proximity. Regional dimension should rather be seen as a merger of institutions under a single administrative leadership (Karelian component – Karelian

Research Center of RAS, Siberian component – Siberian Branch of RAS, Far Eastern component – Far Eastern Branch of RAS, and finally, the national component – Russian Academy of Sciences).

As a result, administrative actions implemented within the framework of regional dimension provide a more prospective explanation of the behavior of the webgraph dynamics model proposed. As stated in Pechnikov (2010), the web resources of research institutions belong to the so-called "regulated web resources". This means there is an official document that sets out the aims and objectives of the web resource, the main structural components, rules for changing information contained in them, etc. Manageability of processes on the web refers to implementation of managerial decisions in the form of regulations defining implementation of such decisions. In this case, managerial decisions can be developed for single sites and their communities with the aim of improving their presence on the Web. Such managerial decisions, for example, may be rules requiring site developers to create hyperlinks pointing to the websites of higher and/or subordinate institutions.

Analysis of the proposed webgraph dynamics model from the point of view of administrative action provides the following results.

At time $t_1$, emergence of Karelian component $A$ dates back to 2007 as evidenced by the creation of independent sites of institutions in the period 2006-2007 (Internet Archive, 2015). The authors are aware that the program for creation of a unified information system for the Karelian Research Center of RAS was launched in 2005. Implementation of the program is what can be considered as an administrative action that led to this result.

Emergence of the Siberian component was dated to time $t_1$, and by time $t_5$, it already consisted of 4 institutions. The program "*Integrated information and telecommunication systems and networks, telecommunication and information resources, information processes in systems and networks*" was launched in 2007 as part of fundamental research at the Siberian Branch of RAS (Program, 2007-2009). This was an administrative action that resulted in the development of the Siberian component (CO PAH, 2015).

Target program of the Presidium of RAS "*Informatization of Research Institutions and the Presidium of RAS*" was launched in 2002. This program represents an administrative action, which, by 2006, led to the creation of such section on the RAS website as "*Information Systems research institutions of the Russian Academy of Sciences*" (http://ras.ru/sciencestructure/informationsystems.aspx), which is actually a directory of the websites of RAS institutions. It remains 2-3 years before a large connected component containing 43 sites emerges. For now, backlinks pointing from websites of the institutions to RAS site have emerged (Figure 5). Thus, time $t_8$ corresponds approximately to 2009. This

statement is also supported by studies (Pechnikov, 2010) conducted on the connectivity of the web resources of RAS institutions at that period of time.

Far Eastern component $R$, containing 6 sites, is detected at time $t_9$ (Figure 6). It emerged thanks to the creation of the section "*Research centers and institutions of the Far Eastern Branch of RAS*" on the website of the Far Eastern Branch of RAS (http://www.febras.ru/instituty.html), whose emergence is dated from the year 2012 (Internet Archive, 2015). In turn, it can be assumed that this section emerged due to such administrative action as constant monitoring of the information space of the Far Eastern Branch of RAS, which is led by one of the vice-chairmen of the Far Eastern Branch of RAS (Khanchuk & Naumova, 2009).

From the above, the following conclusions can be drawn. The main stages in the behavior of the proposed dynamics model of the webgraph of RAS institutions are well supported by well-known administrative actions implemented as an information system development program at the level of regional research centers and branches, and RAS in general. Moreover, a comparison of the sequence of model instants of time $t_1$, $t_2$, ..., $t_{10}$ with actual dates of adoption of information system development programs and their implementation, which are an administrative action on web resources, shows that there is an almost perfect match. Analysis of the model and administrative action enables us to date the main model instants of time with a high degree of reliability: namely, time $t_1$ corresponds to 2006, $t_5$ to 2007, $t_8$ to 2009, $t_9$ to 2012, and $t_{10}$ to 2014. Based on such dating, other less significant model times can be approximated.

## Conclusion

In this paper, we have proposed a dynamics model of the webgraph of the research institutions of the Russian Academy of Sciences. The model is based on fixation of a webgraph built for a given (finite) instant of time and "return back" through removal of multiple edges (hyperlinks). This is done using a webgraph constructed from data collected in 2014 on hyperlinks interlinking 397 websites of institutions. The dynamics model was built based on the assumption that sequential removal of multiple edges connecting any pair of websites allows to build a webgraph that reflects the state of RAS webspace some time ago.

Thus, we get a "time series" of webgraphs, which reflects the development of RAS webspace on a time interval. The time series is provisionally divided into 10 segments. Results stored in archive of websites "*Internet Archive: Wayback Machine*" enables one to establish the real dates of both the model instants of time and creation of the websites of RAS institutions.

Analysis of the constructed model shows that the main development trend of RAS web space features exponential increase in the number of edges of the webgraph for every next time interval. This in turn expands the maximum connected component through absorption of smaller

components and emergence of new small components. The studies conducted have shown that emergence of new links at each time step in the webgraph dynamics model cannot be explained only by the well-known principles of preferential attachment and initial attractiveness of nodes. Much more important are the so-called 'administrative actions', namely the target programs aimed at developing information resources both at the regional and national level. The study showed an almost perfect match between the model dates corresponding to significant changes in the webgraph dynamics model and the real dates when the administrative actions came into force.

The study was limited to one country (Russia) and one type of establishments (research institutions of the Russian Academy of Sciences). However, this research object is a typical example of regulated set of scientific websites – a set that can realistically be subjected to administrative action. Therefore, the results of the study are first and foremost important in the formation of the web space of research institutions in Russia, because the very system of research institutions is today under reform. Equally important is the fact that the results obtained can be useful to heads of governmental authorities in the field of science in developing countries, especially those where the use of web technology as a means of communication is still in the early development stages.

## References

Bollobas, B., & Riordan, O. (2003). Mathematical results on scale-free random graphs. *Handbook of graphs and networks*, Wiley-VCH, Weinheim, pp. 1-34.

Gephi (2015). The Open Graph Viz Platform. Retrieved May 11, 2015, from https://gephi.github.io

Golovin, A.S. & Pechnikov, A.A. (2013). Database of external hyperlinks for web fragments investigation. *XXI Century University Infomedia: Proceedings of the VII International Scientific and Practical Conference*, September 23-27, 2013, *Petrozavodsk*, 2013, pp. 55-57.

Internet Archive (2015). Internet Archive: Wayback Machine. Retrieved May 11, 2015, from https://archive.org/web

Khanchuk, A.I., & Naumova, V.V. (2009). Information space of the Far Eastern Branch of the Russian Academy of Sciences. *Bulletin of the Far Eastern Branch of the Russian Academy of Sciences*. No. 4, pp. 122-129. Retrieved May 11, 2015, from http://www.cnb.dvo.ru/vestnik/index_eng.htm

Pechnikov, A.A., & Chernobrovkin, D.I. (2014). Adaptive crawler for external hyperlinks search and acquisition. *Automation and Remote Control*, 75(3), 587-593.

Pechnikov, A.A. (2010). Research methods for regulated thematic web fragments. *Proceedings of the Institute of Systems Analysis, Russian Academy of Sciences*. Series: *Applied problems of macro systems management*. Vol. 59, pp. 134-145.

Russian Academy of Sciences. (2015). In *Wikipedia, The Free Encyclopedia*. Retrieved May 10, 2015, from https://en.wikipedia.org/wiki/Russian_Academy_of_Sciences

Russian Government Decree (RGD). (2015). *Распоряжение Правительства Российской Федерации от 30 декабря 2013 г. N 2591-р г. Москва* [Order of the Government of the Russian Federation dated December 30, 2013 N 2591-p Moscow]. No. 2591-p of December 30, 2013. Retrieved May 11, 2015, from http://www.rg.ru/2014/01/09/fano-site-dok.html

Webometrics Ranking for Research for Institutions in Russia (WRR) (2015). Retrieved May 11, 2015, from http://webometrics-net.ru

CO PAH (2015). Программа 4.5.1. *Интегрированные информационно-телекоммуникационные системы и сети, телекоммуникационные и информационные ресурсы, информационные процессы в системах и сетях*. [Program 4.5.1. Integrated information and telecommunication systems and networks, telecommunication and information resources, information processes in systems and networks (2007-2009)]. Retrieved May 11, 2015, from http://www.sbras.ru/cmn/onr/proj.php?id=50&inst=1

---

*Bibliographic information of this paper for citing:*

Pechnikov, Andrey A., & Nwohiri, Anthony M. (2015).   "Webgraph connectivity and dynamics: Russian research institutions."   *Webology*, 12(1), Article 135. Available at: http://www.webology.org/2015/v12n1/a135.pdf

---