Toward the Evaluation of Research Groups based on Scientific Co-authorship Networks: The Robocorp Case Study

¹Micael Couceiro; ²Filipe Manuel Clemente; and ³Fernando Martins

¹Institute of Systems, University of Coimbra, Portugal.

²Faculty of Sport Sciences and Physical Education, University of Coimbra, Portugal.

³ Instituto de Telecomunicações, Covilhã, Portugal.

ID # (2714) Received: 24/12/2012 In-revised: 04/02/2012 Corresponding Author: Filipe Manuel Clemente e-mail: filipe.clemente5@gmail.com

KEYWORDS

Co-authorship Networks; Graph Theory; Researchers connectivity; Collective Evaluation.

ABSTRACT

Scientific cooperation is one the most important issues to improve the research quality. A multidisciplinary scientific group connection among different knowledge areas (e.g., engineering, mathematics, sports, sociology and others) can be a potential factor to build skilled manpower necessary for strong scientific research. Therefore, based on a case study from *Robocop*, a multidisciplinary group with researchers from several scientific fields, this paper presents the scientific cooperation between researchers through networking graph theory. These networks are addressed to answer a broad variety of questions about collaboration patterns, such as the number of papers authors write, with how many researchers they write and how researchers "connect" to make papers in specific areas. First, a weighted adjacency matrix is built based on papers published in accordance with international standards (e.g., ISBN, ISSN), in which it is possible to perceive the connectivity among researchers. Secondly, an easy-to-use Mat Lab script was developed to compute the data, thus presenting the scientific networks. Afterwards, in order to further study the sub communities inside the research group, a graph partition methodology was used to divide the graph into clusters. Moreover, several network concepts were used to evaluate the intra and inter-researchers performances as well as the collective performance of the whole group. Results showed that the research group is integrally connected when considering all published papers. However, dividing the networks by scientific areas, one can observe that some researchers 'loses' their connectivity, i.e., some authors only publishes on specific scientific categories or with specific researchers within the group.

تقييم أداء مجموعات مشاريع البحث العِلمي المُؤسس على شَبكات البحث المُشارك: أنموذج معهد روبوكوب Robocorp Institute في جمهورية البرتغال

¹ **مايكل كوكيرو**، ²**فيليب مانويل كليمنت**، و ³ **فيرناندو مارتيز** ¹ معهد النُظم ، جامعة كويمبرا ²كلية علوم الرياضة والتربية البدنية ، جامعة كويمبرا ³ معهد نظم هندسة الاتصالات وتقنية المعلومات ، قسم الهندسة الكهربائية، معهد الهندسة ، كويمبرا، جمهورية البرتغال

المستلخص

رقم المسودة: # (2714) تاريخ إستلام المسودة: 2012/12/24 تاريخ المسودة المُعَدَلة: 2013/02/04 الباحث المُرَاسَل: فيليب مانويل كليمنت بريد إلكتروني: filipe.clemente5@gmail.com

الكلمات الدالة

شبكات بحث مشترك؛ نظرية الرسم البياني؛ الربط بين الباحثين؛ تقييم جماعي.

يعتبر التعاون العلمي من أهم العوامل لتحسين نوعية ومستوى البحوث العلمية، حبث أن اتصال المجموعة العلمية مُتَعددة التخصصات بين مجالات المعارف المختلفة ،على سبيل المثال، الهندسة والرياضيات والرياضة وعلم الاجتماع، وغيرها، يمكن أن يكون عاملا أساسياً لبناء وتطوير المقدرات الماهرة اللازمة لإعداد مشروعات البحث العلمي الرصين. لذلك، واستنادا إلى در اسة حالة من معهد رويوكوب Robocorp Institute، وهي مجموعة متعددة التخصصات مع باحثين من عدة مجالات علمية، تستعرض هذه الورقة التعاون العلمي بين الباحثين من خلال (نظرية شبكات الرسم البياني Networking Graph Theory). تعالج هذه الشبكات مجموعة واسعة من الأسئلة حول أنماط التعاون، مثل عدد صفحات وعاء البحث العلمي ، وعدد الباحثين المشاركين وكيفية التواصل بين هؤلاء الباحثين المشاركين ومستويات مشاركاتهم ولتحقيق هكذا التواصل يتم أولا بناء مصفوفة التجاور المرجح وتقييم جدوي وأهمية المشاركة على أساس مستويات الأبحاث المنشورة للباحث المعنى وفقا للمعايير الدولية، والتي منها مستوي مصادر النشر (Evaluation of Resources). أما ثانيا يتم تطوير وسيلة سهلة الاستخدام مختبر حساب البيانات (Mat Lab)، لتكملة البيانات وبالتالي تكملة تواصل الشبكات البحثية العلمية. بعد ذلك، ولأغراض مواصلة دراسة المجتمعات الفرعية داخل مجموعة البحث، يتم استخدام منهج تقسيمات الرسم البياني Methodology) of Graph Partition) القائم على مجموعات بيانات البحث العلمي لتقسيم الرسم البياني إلى مجموعات. إضافةً إلى ذلك، تم استخدام العديد من المفاهيم شبكة لتقييم الأداء العام وفيما بين الباحثين، فضلا عن الأداء الجماعي للفريق بأكمله. وأظهرت النتائج أن كل مجموعة فريق بحثي يرتبط ويتوافق عند النظر في جميع الأوراق المنشورة لهم ومع ذلك عند تقسيم الشبكة إلى مجالات ومحاور علمية يمكن ملاحظة فقدان بعض الباحثين لميزة التواصل بينهم ، وتقسيم الشبكات عن طريق المجالات العلمية، يمكن للمرء أن يلاحظ أن بعض الباحثين يفتقد ميزة التواصل مما حصر مخرجاتهم البحثية في محاور محددة ومشاركاتهم مع باحثين محددين.

(1) Introduction

Graph theory was created through Euler's solution of the Königsberg's bridges problem in 1736. Over time, they have also become extremely useful as representation of a wide variety of systems in different scientific areas. Biological, social, technological, and information networks can be studied as graphs, and graph analysis has become crucial to understand the features of these systems (Fortunato and Castellano, 2009). One such property is the "small world effect", which is the name given to the finding that the average distance between vertices in a network is short usually scaling logarithmically with the total number n of vertices (Girvan and Newman, 2002).

(1.1) Co-authorship Network

Despite of different applications of network theory, very little attention has been given in the evaluation of research groups. The coauthor ship of a paper can be understood as a documenting collaboration between two or more authors, in which these collaborations form a co-authorship network. According to Newman (2004b) scientific publications can be better represented by social networks than many affiliation networks. This may imply that researchers who have written more papers together are genuinely acquainted with themselves. Therefore, the study of co-authorship networks has received some of the attention of researchers, seeking to understand the dynamics inherent to the scientific activities and collective work (e.g., Okubo, Miquel, Frigoletto and Doré, 1992; Newman, 2004a, 2004b; Cardillo, Scellato and Latora, 2006).

Some studies had used scientific databases in order to understand the connectivity and dynamics among researchers and among scientific areas (*e.g.*, Narin, Stevens and Whitlow, 1991; Arunachalam, Srinivasan and Raman, 1994; Newman, 2004a, 2004b; Rodriguez and Pepe, 2008). Over time, some studies analyzed the collaboration among researchers through co-authorship networks (Yoshikane and Kageura, 2004).

Several different methodologies to measure the connectivity among authors have been presented. For instance, some studies are based on the strength between link nodes, using these indicators to measure the co-authorship networks (Narin, 1991; Arunachalam et al., Kundra and Kretschmer, 1999). However the reported studies were mainly based on static network and did not analyzed the evolution of networks over years or over the increasing number of published scientific articles. In response to this, Yoshikane and Kageura (2004) used the Monte-Carlo simulation to evaluate the growth and change of networks. Nevertheless, the authors have not examined the observed accumulation according to time series, *i.e.*, their analysis was based on simulation data instead of real data. Similarly, Newman (2004a, 2004b) studies showed that the structure of such networks turns out to reveal many interesting features between academic communities. Their studies analyzed results associated to the number of authors, number of papers per author, number of authors per paper and clustering coefficients. Additionally, in order to measure the complex patterns, the author used several graph theory techniques such as the shortest paths, average distances and the weighted collaboration networks.

(1.2) Scientific Contribution

However, the majority of the works published about co-authorship networks focused the study on scientific databases, not considering the monitoring and evaluation of research centers/groups. As a result, this paper presents an evaluation strategy based on graph theory to further evaluate research groups. Besides presenting the connectivity between researchers as Newman's work (2004), the herein proposed methodology also offers a graphical way to identify the most contributing researchers (*i.e.*, the ones that publishes the most within the group) and the group partitioning (*i.e.*, identification of possible subcommunities within the group).

To that end, an easy-to-use script developed in *MatLab* allows the user to add new papers, thus updating the network graphs of the research group and possible subcommunities within the group. Experimental results were obtained by using data from a multidisciplinary group with researchers from several scientific fields named as *RoboCorp*. Only the last five years were analyzed, from 2007 to 2011, with a total of 108 original papers published in accordance with international standards (*e.g.*, ISBN, ISSN). Also, only twelve researchers were considered since they were the only ones publishing since the beginning of 2007 until the end of 2011.

(2) Scientific Cooperation

In order to achieve the proposed objectives previously defined, this work may be divided into three research aspects: *i*) the graphical representation of co-authorship networks based on researchers' interactions that arises from the publication of scientific articles from a given research group; *ii*) the network partitioning of the co-authorship network into multiple sub-groups based on researchers' connectivity; and *iii*) the evaluation of the co-authorship network through network concepts (*i.e.*, network indices) to describe the topological properties of the research group.

(2.1) Co-authorship Networks

The concept of co-authorship networks was first introduced by Newman (2004) in which network theory was used to represent the scientific collaboration between researchers. The coauthorship networks herein presented are similar to the ones introduced by Newman (2004) but contributes with the development of an easy-touse script that allows to graphically represent the relation between researchers and further identify the most contributing researchers within a research group of n researchers (*i.e.*, the ones that publishes the most within the group).

A *MatLab* script denoted as *wgPlot* was developed by Michael Wu (2009) which allowed to plot graphs similarly to *gPlot*, a *MatLab* function that allowed to plot *n* nodes connected by links representing a given *adjacency matrix* $A = [a_{ij}] \in \mathbb{R}^{n \times n}$ defined by:

$$a_{ij} = \begin{cases} 1, connection between node i and j \\ 0, no connection between node i and j \end{cases}$$
(1)

Within co-authorship networks context, edges are adjacent when there is one vertex incident

with them, *i.e.*, the connection between vertices (*i.e.*, researchers) is defined by the co-authorship of a single paper published in accordance with international standards (*e.g.*, *ISBN*, *ISSN*). It is noteworthy that in our special situation, in which each adjacency matrix represent a publication, the diagonal elements (*i.e.*, when i = j) are set equal to 1 to identify the researcher as one of the coauthors of the publication. As an example, consider the herein presented publication in which the first author corresponds to the first vertex and so on. The researchers, (*i.e.*, n = 5), but the two last did not contribute to this work. The adjacency matrix of this publication would be represented by:

$$A = \begin{bmatrix} 11100\\ 11100\\ 11100\\ 00000\\ 00000 \end{bmatrix} \in \Re^{5 \times 5}$$
(2)

The script *wgPlot* from Michael Wu (2009) allows the user to input an adjacency matrix with weighted edges and/or weighted vertices being denoted as *edge-weighted edge-adjacency matrix* A_w , introduced by Ernesto (1995).

The weighted matrix A_w can be easily defined by the sum of all adjacency graphs each one generated by a single publication. Let us suppose an example in which two other more publications, besides this herein presented. The first, second, fourth and fifth researchers handle the co-authorship of the first publication while the second publication has the contribution from the first, second and fifth researchers. The matrix $A_w = [w_{ij}] \in \Re^{5\times 5}$ would then be represented by:

$$A_{w} = \begin{bmatrix} 33112\\ 33112\\ 11100\\ 11011\\ 22012 \end{bmatrix}$$
(3)

To allow a graphical representation of the scientific cooperation, the script presented by Michael Wu (2009), denoted as *wgPlot*, was further extended based with the following features:

- (i) The edge (*i.e.*, researcher) size i, i = j, of the network is proportional to the number of publications in which he/she is a coauthor w_{ii}.
- (ii)The vertex (*i.e.*, cooperation between researchers) thickness and colormap of the network is proportional to the number of publications in which researcher i and j, $i \neq j$, publish together.
- (iii) The script receives as input a binary database (*e.g.*, excel file) in which each line corresponds to a publication and each column to a researcher, *i.e.*, each line corresponds to an adjacency matrix A.
- (iv) Besides returning the network from A_w , it also returns the clusters, *i.e.*, subcommunities, of the research group based on Hespana's work (Hespana, 2004) and extensively used in (Lim *et al.*, 2005). This last point will be further explained in next section.

(2.2) Network Partitioning

In order to detect groups among researchers, graph theory has specific methodologies to constitute partitions. Uniform graph partition consists of dividing a graph into components, such that the components are of about the same size and there are few connections between the components. One of the functionalities of the graph partition is generate communities. Communities, also called clusters or modules, are groups of vertices which probably share common properties and/or play similar roles within the graph (Fortunato, 2010).

The uniform graph partition has gained importance due to its application for clustering and detection of cliques in social, pathological or biological networks (Fiduccia and Mattheyses, 1982). Commonly the graph partition is defined by G = (V, E) where V is the vertex and E is the edges, such that is possible to partition G into smaller components with specific properties. A k-partition of V is a collection $P = \{V_1, V_2, ..., V_k\}$ of k disjoint subsets of V, whose union equals V (Hespanha, 2004). The *MatLab* function *grPartition* described in the technical report of Hespana (2004) allows the fast partition of large graphs. This function implements a graph partitioning algorithm based on spectral factorization. The herein proposed *MatLab* script then merges the *wgPlot* and *grPartition* functions, with a few adaptations as previously presented, to understand the scientific cooperation patterns within a given research group, such as the numbers of papers authors write, how many colleagues they write them with and the existences of sub communities among them.

Therefore, running the script with the previously described example (*cf.*, section Co-authorship Networks) would then return the following co-authorship network, thus identifying the scientific cooperation among researchers (Figure 1).



Figure 1. Co-authorship Network between 3 Publications from 5 Researchers (example previously described and represented by A_w in equation 3).

As one can observe, the first and second researchers are the ones that publishes the most (*i.e.*, larger vertices). This may be closely related with the cooperation between both researchers since they present a high connectivity (*i.e.*, ticker edge). On the other hand, the third and fourth researchers are the ones that publishes less (*i.e.*, small vertices). Also, there are essentially two subgroups (*i.e.*, vertices of different colors) in which one is formed by the third researchers and the other by the remaining researchers, *i.e.*, first, second, fourth and fifth researchers. This could mean that

the third researcher is not cooperating as other group members are.

(2.3) Network Concepts

Many kinds of networks (*e.g.*, biological, sociological or others) share topological properties. To identify and describe such properties most potentially useful network concepts are known from graph theory. In co-authorship networks context, one can divide network concepts into:

- (i) Intra-researcher network concepts (*i.e.*, network properties of a node).
- (ii) Inter-researcher network concepts (*i.e.*, network relationship between two or more vertices).
- (iii) Group network concepts (*i.e.*, whole network concepts).

To allow using most of the network concepts, one can create a new relative weighted adjacency matrix, $A_r = \begin{bmatrix} r_{ij} \end{bmatrix} \in \Re^{n \times n}$ defined as:

$$r_{ij} = \begin{cases} \frac{w_{ij}}{\max A_w}, i \neq j \\ w_{ij}, i = j \end{cases}$$
(4)

where $0 \le r_{ij} \le 1 \le 1$ for $i \ne j$, with i, j = 1, ..., n

The denominator $\max_{i \neq j} A_w$ corresponds to the larger connectivity between researchers, *i.e.*, the researchers that most published together.

It is noteworthy that the diagonals of A_r will still represent the number of papers published by a researcher. However, this value is not considered to compute the network concepts herein presented. (2.3.1) Intra-Researcher Concepts

The first concept and one of the widely used in the literature for distinguishing a vertex of a network (*cf.*, (Horvath, 2011), is the *connectivity* (also known as degree). The connectivity of researcher i can be defined by:

$$k_i = \sum_{i \neq j} r_{ij} \tag{5}$$

such that $k = [k_i] \in \Re^{1 \times n}$ is the vector of the connectivity of researchers.

In the situation herein presented, *i.e.*, coauthorship networks, the connectivity equal the sum of connection weights between researcher iand the other researchers. The most cooperative researcher, or researchers, can be found by finding the index/indices of the maximum connectivity.

$$k_{\max} = \max_{i} k_{j} \tag{6}$$

Therefore, one can define a relative connectivity, known as *scaled connectivity*, of researcher i as:

$$S_i = \frac{k_i}{k_{\max}}$$
(7)

such that $s = [s_i] \in \Re^{1 \times n}$ is the vector of the relative connectivity of researchers.

In research group context, one could interpret the scaled connectivity as a measure of cooperation level of a given researcher in which high values of s_i (*i.e.*, as s_i tends to 1) indicate that the i^n researcher works with most of the other researchers from the group. However, a researcher may present a high connectivity but may be unable to produce consensus among his/her coauthors. In other words, he/she may publish with several other researchers that do not publish with each other. Therefore, the *clustering coefficient* of researcher *i* offers a measure of the degree of interconnectivity in the neighborhood of researcher *i*, being defined as:

$$C_{i} = \frac{\sum_{j \neq i} \sum_{j \neq i} r_{ij} r_{jl} r_{ki}}{\left(\sum_{j \neq i} r_{ij}\right)^{2} - \sum_{j \neq i} (r_{ij})^{2}}$$
(8)

such that $c = [c_i] \in \Re^{1 \times n}$ is the vector of the clustering coefficient of researchers.

As Watts and Strogatz (1998) suggests, this intra-researcher network concept is a density measure of local connections, or "cliquishness". Hence, the higher the clustering coefficient of a researcher, the higher is the scientific cooperation among its coauthors. In other words, a clustering coefficient tends to zero if all the coauthors of a given researcher do not publish much with each other.

The relationship between the clustering coefficient and the connectivity has been used to describe structural (hierarchical) properties of networks *e.g.*, (Ravasz *et al.*, 2002). Despite that in most situations the clustering coefficient is inversely related to the connectivity, researchers that are associated with only one scientific category

may have a high connectivity and also a high clustering coefficient. However, researchers that work in multiple scientific categories are correlated with a larger number of researchers, but many of these researchers do not publish with each other, leading to a smaller clustering coefficient.

As a multidisciplinary research group, a weighting distribution of the cluster coefficient and the connectivity between researchers should be taken into account. RoboCorp research group considers that it is important to have a high level of connectivity since researchers should present a persistent partnership over time, seeking to create strong and lasting relationships. At the same time, it also considers that the cluster coefficient of a researcher is relevant to the group since it is necessary to produce partnerships in order to create interdisciplinary relationships, thus increasing the collective productivity of the group, *i.e.*, the group should sustain itself as a group avoiding to ensure exclusive priority to the individual performance. Therefore, a weighting function, denoted as global rank, was defined as:

$$g_i = \rho_s s_i + \rho_c c_i \tag{9}$$

where $\rho_s + \rho_c = 1$, such that $g = [g_i] \in \Re^{1 \times n}$ is the vector of the global rank of researchers.

Note that the scaled connectivity ^Si was chosen over the unscaled one k_i since it lies between 0 and 1 as the clustering coefficient, thus resulting in $0 \le g_i \le 1$. Taking into account that the main objective of RoboCorp research group is to give priority to the collective performance i.e., (overall interaction between researchers), one can ponder a balanced consideration of $\rho_k = \rho_c = 0.5$. The topranked researcher, i.e., the one presenting the higher g_i , will then be denoted as the researcher centroid. Within research group context, the researcher centroid could be considered as a hierarchically superior member (e.g., supervisor).

As a result, the herein proposed script returns the scaled connectivity, clustering coefficient and global rank of researchers of a given co-authorship network. Using the previously presented example, the script would return the following output:

<i>s</i> =	[1.0000	1.0000	0.2857	0.4286	0.7143]
<i>c</i> =	[0.3922	0.3922	1.0000	0.7778	0.6667]
<i>g</i> =	0.6961	0.6961	0.6429	0.6032	0.6905]

As it is possible to see, researcher 3 presents the lower connectivity of the group and the higher clustering coefficient (since researchers 1 and 2 highly cooperates with each other), thus resulting in a global rank higher than researcher 4 (which present an higher connectivity than researcher 3). It is also noteworthy that in this specific example, both researchers 1 and 2 are considered the researchers centroids. Within co-authorship context, they could both be considered as hierarchically superior to the other researchers (*e.g.*, supervisors).

(2.3.2) Inter-Researcher Concepts

To complement the intra-researcher concepts, at least two inter-researcher concepts need to be considered. The first one arises from the researcher centroid (defined above) in which one can express his/her connection strength to all other researchers as:

$$CC_{i,centroid} = \begin{cases} r_{i,centroid}, i \neq j \\ 1, i = j \end{cases}$$
(10)

This inter-researcher concept is denoted as *centroid conformity* and corresponds to the adjacency between the researcher centroid and the i^{th} researcher, such that $cc = [cc_i] \in \mathbb{R}^{1 \times n}$ is the vector of the centroid conformity of researchers. In other words $^{CC_{i,centroid}}$, presents the cooperation level of the i^{th} researcher with the top-ranked researcher.

The second inter-researcher concept is based on the topological overlap presented in several works such as (Ravasz *et al.* 2002) and (Horvath, 2011) which represents the pair of researchers that cooperates with the same researchers. However, this measure presents the overlap between two researchers even if they do not publish with one another. In other words, the topological overlap between the i^{th} researcher and the j^{th} researcher depends on the number of published papers with the "shared" researchers but does not take into account the number of published papers between them as it should. Moreover, the topological overlap is represented by a symmetric matrix, thus presenting

the overlap between researchers but neglecting the most independent researcher of the pair. Therefore, by using the concepts inherent to the clustering coefficient (see, equation 7), one should consider not only the "shared" papers but also the influence of the conjoint publications among researchers i and *j*. In other words, if two researchers publish with the same other researchers, then the cooperation between both of them allows building triangular relations between the other researchers. However, the i^{th} researcher may be more dependable from the *i*th researcher if he/she only publishes with the same researchers than researcher jth which, in turn, is able to publish with other researchers. As a result, similarly to Ravasz et al. (2002) and Horvath (2011), one can define a topological dependency $T_d = [td_{ij}] \in \mathbb{R}^{n \times n}$ as:

$$td_{ij} = \begin{cases} \frac{\sum_{l \neq i, j} r_{il} r_{lj} r_{ij}}{\sum_{l \neq i} r_{il}}, i < j \\ \frac{\sum_{l \neq i, j} r_{il} r_{lj} r_{ij}}{\sum_{l \neq i} r_{lj}}, i > j \\ 1, i = j \end{cases}$$
(11)

with i, j, l = 1, 2, ..., n.

As a consequence, two researchers have a high topological dependency, *i.e.*, $td_{ij} = 1$, if they publish with the same researchers and with one another. In other words, the more researchers are "shared" between two researchers that highly publish with one another, the stronger are their cooperation and more likely they will both represent a small cluster.

However, since T_d corresponds to a square matrix with the size equal the number of researchers and since that contrarily to the adjacency matrix or topological overlap usually used in the literature, *e.g.*, (Horvath, 2011), T_d is not symmetric, *i.e.*, $td_{ij} \neq td_{ji}$, it makes it difficult to compare the td_{ij} and td_{ji} pais. Therefore, one can introduce a new inter-researcher concept denoted as *topological inter-dependency* $T_{id} = [ti_{ij}] \in \mathbb{R}^{n \times n}$ as:

$$T_{id} = T_d - T_d^t \tag{12}$$

Wherein T_d^T is the transpose of matrix $T_d \& T_{di}$ corresponds to an ant symmetric square matrix, *i.e.*, $ti_{ii} = -ti_{ii}$.

In co-authorship networks, one can easily observe dependencies between researchers such that if $t_j > 0$ then the i^{h} researcher depends on the j^{h} researcher to publish with his/her coauthors. Moreover, when associated to other network concepts, *e.g.*, (researcher centroid) the relative topological dependency allows identifying possible dependencies between researchers and even hierarchical relations.

As a result, the herein proposed script returns the centroid conformity as well as the topological overlap of a given co-authorship network. Using the previously presented example, the script would return the following output:

 $cc = \begin{bmatrix} 1.0000 & 1.0000 & 0.3333 & 0.3333 & 0.6667 \end{bmatrix}$

	0	0	- 0.1190	-0.1058	- 0.0889]
	0	0	- 0.1190	-0.1058	- 0.0889
$T_{id} =$	0.1190	0.1190	0	0	0
	0.1058	0.1058	0	0	0.0593
	0.0889	0.0889	0	- 0.0593	0

In this example, one can observe that the most cooperative researchers (researchers 1 and 2), *i.e.*, the one that presents the higher connectivity, are also the ones that most cooperate with each other (as they are both the researchers centroids). Although this is not a linear relationship, it is highly probable that when the group has more than one researcher centroid, they highly cooperate with one another. It is also possible to highlight researcher 5 as it seems to cooperate more with researchers centroids. This could mean that he/ she is hierarchically superior to researchers 3 and 4, *i.e.*, (co-supervisor). Through the topological inter-dependency measure, one can, for instance, observe that researcher 3 presents a high overlap with researchers 1 and 2 - it could be concluded that approximately 10% of its publications highly depends on them. Looking at Figure1 it is easy to observe that while researcher 4 only publishes with the same coauthors of researchers 1 and 2, they are both able to publish with other researchers, *i.e.*, (researcher 3). Despite considering a small database, this may induce a "master-slave" relation between researchers 1 and 2 and researchers 4. In other words, researchers 1 and 2 could be considered hierarchically superiors to researcher 4 (*e.g.*, supervisor-student relationship).

(2.3.3) Group Concepts

Although both inter and intra-researcher concepts are useful to identify properties between researchers, group network concepts also need to be considered to achieve properties of the full research group.

The inter-researcher connectivity allows retrieving several other group network concepts such as the *network density* which can be defined as:

$$D = \sqrt{\frac{\sum ki}{n(n-1)^2}}$$
(13)

Within co-authorship networks, the density measures the overall cooperation among researchers. A density that tends to 1 indicates that all researchers strongly publish with each other.

Another network concept based on the connectivity of researchers is the *network heterogeneity* which is closely related to the variation of connectivity across researchers (*cf.*, (Albert, Jeong and Barabasi, 2000) and (Watts, 2002). As Horvath's work (2011), it is herein defined as the coefficient of variation of the connectivity distribution:

$$H = \sqrt{\frac{n\sum k_i^2 - (\sum k_i)^2}{(\sum k_i)^2}}$$
(14)

Since the heterogeneity measure is invariant with respect to multiplying the connectivity by a scalar, one could use the scaled connectivity instead of the connectivity. Many complex networks have been found to exhibit an approximate scalefree topology, which implies that these networks are very heterogeneous. In other words, a high heterogeneity of the co-authorship network means that the research group exhibits a high level of subcommunities and there is, collectively, a low level of cooperation between researchers.

Finally, to further analyze the co-authorship network, a widely used measure denoted as *network centralization* was used. The network centrality (or degree centralization as Freeman, 1978) addresses) can be defined as:

$$C = \frac{n}{n-2} \left(\frac{\max k}{n-1} - D \right)$$
(15)

A centralization of the co-authorship network close to 1 means that one researcher strongly cooperates with all other researchers which, in turn, present a small (or inexistent) cooperation with each other. In contrast, a centralization of **0** indicates that all researchers cooperates equally between each other.

As a result, the herein proposed script returns the network heterogeneity, density and centralization of a given co-authorship network. Using the previously presented example, the script would return the following output:

$$H = 0.4751 \quad D = 0.4000 \quad C = 0.3056$$

In brief, one could conclude that the group is more or less homogeneous. For instance, researcher 3 should cooperate more in order to reduce the heterogeneity of the group. Yet, the group presents a low density since researchers do not cooperate enough with each other. Although the group presents a low network centralization, one need to be aware that only three publications from five researchers were analyzed, thus making it impossible to observe a large discrepancy between them. Therefore, to further endorse the current methodology, next section presents an extended analysis of *RoboCorp* group, benefiting from the full properties of the herein proposed script.

(3) Case Study

The previously defined methodology was applied to the data obtained from a multidisciplinary group with researchers from several scientific fields named as *RoboCorp*. Only the last five years were analyzed, from 2007 to 2011, with a total of

158 original papers published in accordance with international standards (e.g., ISBN, ISSN). Also, only twelve researchers were considered since they were the only ones publishing since the beginning of 2007 until the end of 2011. It is also important to emphasize that although the group presents multiple supervisor-student relationships, there is not a pre-defined hierarchy within the research group.

Based on (Harzing, 2007), RoboCorp fields of studies can be divided into the following two categories:

(1) Engineering, Computer Science and Mathematics.

(2) Social Sciences, Arts and Humanities.

Therefore, Table 1 depicts a summary of *RoboCorp* publications in the last 5 years.

Categories, Fields of Studies	Category (1)	Category (2)	Global
Authors / Category	06	06	012
Number of Papers	78	30	108
Papers / Author	14.25	5.58	09.92
Authors / Paper	02.19	2.23	02.20

Table 1: Summary of RoboCorp Data during the Last 5 years, 2007-2011

Category (1) Engineering, Computer Science & Mathematics

Category (2) Social Sciences, Arts & Humanities

Authors/ Category = number of authors from a category. Number of Papers = total number of papers published in a category.

Papers/Author = mean number of papers published by an author in a category. Authors/Paper = mean number of authors on a paper in a category).

As one can observe, the number of authors in both scientific categories is the same. Also, more than 70% of the published papers fall into category 1 which may infer two assumptions:

- (i) Authors in category 1 publish more than authors in category 2.
- (ii) Authors in category 2 publish more in category 1 than authors in category1 publish in category2.

The same conclusions can be withdrawn by analyzing the number of papers per author. At least, although the number of authors per paper is approximately the same in both categories and fairly small, it is noteworthy those only RoboCorp researchers are considered. In other words, if an author publish a paper with several external authors (i.e., not from RoboCorp), it will be considered that this paper was only published by one author.

However, all this statistical data may hide a large amount of information. For instance, one should note that most of researchers from a specific area may publish in another, e.g., one of the researchers from sports sciences (category 2) published 10 papers in sports science and 2 papers in engineering (category 1). Also, it is not clear what may be the contribution of each researcher to the collective objective, e.g., while one of the researchers published only one paper in the last 5 years, another one published 47 papers. As another example, the statistical data do not show how researchers within the group cooperate with each other as time goes by. In fact, a given researcher may publish a large amount of papers without any other RoboCorp members, thus increasing the number of papers published by the group. However, this may mean that the group is fragmented and the loss of this researcher would imply a major breakdown of the collective performance of the group. Therefore, the scientific cooperation between RoboCorp members will be further analyzed through networking graph theory.

(3.1) Co-authorship Network Evaluation

The evaluation of a research group should meet its evolution over time. In fact, this aspect is intrinsically related to the qualitative analysis of the group. Indeed, only a cumulative analysis can provide the data that allow us to interpret the performance of the group and its evolutionary trend, defining the growing level of the cooperation between researchers. Therefore, the analysis should not be static, but rather a dynamic and evolving assessment standard as Yoshinkan Kageura (2004) present. To that end, the analysis of the last 5 years, 2007-2011, of the research group *RoboCorp* was profiled in order to interpret its evolution over a period of time, thus allowing understanding the

relationships and connectivity between researchers, as well as the usefulness of them (Figure 5).



Figure 2: Evolution of RoboCorp Co-authorship Network during the Last 5 years, 2007 to 2011.

In 2007, which is the date of the group's foundation, most of researchers were students fulfilling their academic degree (BSc/MSc/PhD), hence not devoted solely to research and investigation. It should be noted that in 2007 only four members (researchers 1, 5, 6 and 10) had recently completed their advanced academic progression (PhD). Later this year, there is a constitution of a sub-group (*i.e.*, vertices of different colors) related to the (BSc) thesis project in category 1, Engineering, Computer Science and

Mathematics, established between the supervisor (researcher 1) and three students (researchers 2, 4 and 9). In 2008, a new sub-group is formed which is related with the (MSc) thesis from category 2, Social Sciences, Arts and Humanities, by researcher 3 under the guidance of researcher 10. It is noteworthy that the four researchers who began the scientific production in 2007, *i.e.*, those who finished their (PhD), have not established much cooperative work with other group members. This may be explained by the publications that arise

from their (PhD) work in which other RoboCorp members did not contribute. Thus, we are witnessing an enlargement of the network vertex without any connectivity established between peers. While in 2008 the maximum connectivity was maintained between researchers 2 and 4, in 2009 we are witnessing a change in the higher level of connectivity of the group, justified by the end of the (MSc) thesis of researcher 2 under the guidance of researcher 4. By analyzing the year 2009, there appears an interaction between the two sub-groups previously defined (*i.e.*, a group from category 1, Engineering, Computer Science and Mathematics, and another from category 2, Social Sciences, Arts and Humanities.) created by the cooperative work among peers. It also appears that the remaining network vertices still produce scientifically but do not cooperate among them. The year of change in the interaction between researchers was in 2010. The major part of the group were graduated students, thus resulting in an enlargement of network vertices (*i.e.*, the scientific productivity of researchers), as well as the connectivity between researchers. There is also the maintenance of cluster formed by researchers from category 1. Engineering, Computer Science and Mathematics, and category 2, Social Sciences, Arts and Humanities, with the addition of an element from

category 2 to this group (researcher 11). The lack of a full connected network is due to the fact that researcher 8 is still in his/her graduation process. From the analysis of 2011, there is a sub-division of the group that was initially formed in 2009 as a result of contributions from authors of different categories, reinforcing its specific intervention. It should be noted that the researcher who had joined the group in 2010 (researcher 11) has been isolated in 2011. There has been equally the addition of a new researcher (researcher 12) who started the scientific production with the other members in 2010. The change between group members may be explained by the fact that researcher 12 have cooperated more with researchers 3 and 10 in 2011 than researcher 11.

By analyzing the global network it is not possible to verify the effectiveness of all individual contribution that researchers provides to the group. Thus, applying the intra-researcher global rank (see, Figure 3) one can analyze the trade-off between the cooperation level of a given researcher and the consensus that he/she can create with its co-authors. Thus, as the number of collective contributions arises, one can observe a descending tendency in the researchers that started publishing in 2007.



Figure 3: Global rank of *RoboCorp* researchers from 2007 to 2011.

Table 2 was obtained by analyzing the final network (cumulative network between 2007 and 2011), thus describing the ranking for both scientific categories previously defined.

Category (1), Engineering, Computer Science and Mathematics	Category (2), Social Sciences, Arts and Humanities	Global
Researcher 2	Researcher 10	Researcher 2
Researcher 1	Researcher 3	Researcher 1
Researcher 4	Researcher 12	Researcher 4
Researcher 9	Researcher 2	Researcher 9
Researcher 3	Researcher 5	Researcher 3
Researcher 10	Researcher 9	Researcher 10
Researcher 8	Researcher 4	Researcher 8
Researcher 7	Researcher 11	Researcher 12
Researcher 5	Researcher 1	Researcher 11
Researcher 12	Researcher 6	Researcher 5
Researcher 6	Researcher 7	Researcher 7
Researcher 11	Researcher 8	Researcher 6

Table 2: Global Rank g_i of RoboCorp Researchers.

(Researchers highlighted in blue are from Category (1), Engineering, Computer Science and Mathematics, while researchers highlighted in green are from Category (2), Social Sciences, Arts and Humanities).

This classification is based on the global rank metric ordered in descending order. Theoretically, researchers from a certain category should present a higher rank in their category. Although this assumption is confirmed in the first ranking positions, the classification of a certain category is influenced by the classification of researchers in the other category. In general, category 1, Engineering, Computer Science and Mathematics, presents a higher rank than category 2, Social Sciences, Arts and Humanities, the group presents a greater focus toward category 1 as the global rank is highly influenced by category 1.

In order to present the inter-researchers metrics, a hierarchy graph was defined based on the centroid conformity and directional edges points toward the direction of the researcher for which a given researcher depends based on the topological interdependency (Figure 4). To categorically group the researchers, a uniform distribution was carried out for each quartile. For instance, for a researchers' centroid conformity between 0% and 25% of the distribution (*i.e.*, first quartile), they are placed in the lower level.





(The researcher centroid is identified as the top vertex; the hierarchy is represented by a uniform distribution of the most conforming researchers to the researcher centroid; and the direction of arrows represents the dependency between researchers). One can observe that there is a relationship between the hierarchically location of researchers and the dimension of the vertices (*i.e.*, number of publications), thus presenting a relationship between researchers who publish the most with the centroid and the overall contribution to the group. Some exceptions may be identified, such as researcher 5 in which, even being hierarchically inferior to researcher 3, he/she presents a greater contribution to the group.

It is also easy to conclude that researchers in a lower hierarchically position present a greater dependence over hierarchically superior researchers. However, an exception to this trend can be observed through the analysis of category 2, Social Sciences, Arts and Humanities, wherein researcher 12, although superior to researcher 2, presents a dependency relationship with him/her. This may be justified by the fact that researchers 2 and 12 published a large number of articles with researchers 3 and 10. However, researcher 2 publishes with a larger diversity of other researchers. It is also possible to observe neutral relationships within the lower classified researchers (*i.e.*, two-way edges), *i.e.*, there is no dependency between researchers despite publishing together. Theoretically, researchers with higher academic degree should be the researchers' centroid. Nevertheless, in this case study, the centroid researcher (researcher 2) is a PhD student oriented by researcher 1 as can be seen by the thicker edge between them. However, as researcher 2 is engaged in other works within the group, being an active member in their execution, his/hers position in the group grows faster than the two senior researchers with higher academic degree (researchers 2 and 10). The importance of evaluating the performance of a research group, more than just analyzing the individual performance and the performance among researchers, one needs to analyze the collective performance (Figure 5).



Figure 5: Group Concepts of RoboCorp Co-authorship Network from 2007 to 2011.

This evaluation will depend on the cumulative results over the last five years of the research group. Figure 4 allows a better understanding about the relationship between researchers, being possible to observe that researchers that published with the researcher centroid are the top ranked ones. However, researchers' position, in some specific cases, do not match their position as senior members, wherein the global hierarchy shows that researcher 10 (*i.e.*, supervisor from category 2) is at an inferior level. In other words, although researcher 10 is the research centroid from category 2, his/hers work is too focused on his/ hers field of research (*i.e.*, category 2). Hence, in a global viewpoint of the group, his/hers position substantially decreases since he/she presents a low general engagement.

The network density allows showing a gradually increasing trend showing a greater collective connectivity between researchers, *i.e.*, there is a greater level of cooperation over time. However, this level of cooperation may not be uniform among researchers, *i.e.*, there may be some researchers who cooperatively publish much more than others triggering an increased level of connectivity. Hence, the network heterogeneity was used to assess the diversity of researchers' intervention. Through the analysis of the heterogeneity, it is possible to verify a downward trend resulting from a greater distributed cooperation among researchers, *i.e.*, the connectivity of researchers is gradually more uniform over time. However, there is an exception to the gradual decrease of the heterogeneity from 2008 to 2009. This may be explained by the increasing number of publications from the sub-groups previously formed, by the increase in the network connectivity between researchers who started co-authorship works in 2009 and the lack of research publications from other researchers still under graduation. There is also the centralization of the group over the years that have been increasing gradually. In fact, it was initially small since most of researchers did not publish collaboratively. Over the years, the triggering for a greater diversified cooperation among researchers unleashed centralization over a particular group or researcher.

In general, it appears that the methodologies applied and the associated metrics allows assessing the collective dynamics of the research group RoboCorp. One can verify that the group focuses its production on the scientific category 1, denoting its growing interest in category 2 over the years. Also it is clear that the group tends to create a greater connectivity and centralization (i.e., hierarchy effect) over the years due to the extending of inter-researchers relations. One can also observe a high cooperation level over the years triggering a decrease in the heterogeneity of the group. Regarding the cooperation between researchers, multiple dependency relationships can be observed between researchers. It is interesting to analyze that this dependence, in general, exists in the hierarchically inferior researchers to the higher ranked ones. Additionally, neutral relationships can be more often observed between hierarchical inferior researchers. By applying the intraresearcher global rank analysis, it was possible to observe the trade-off between the cooperation level and the consensus that a given researcher is able to generates with its co-authors. Given the larger number of collective contributions arising from the higher scientific cooperative productivity, a downward trend of the global rank can be observed in researchers who started such publications in 2007 instead of those who started the publication process later.

In summary, this evaluation methodology allows a global overview of the research group. First, it is possible to observe the personal contribute to the collective objective of the group. By analyzing the global rank, it was possible to see that researchers 5, 6 and 7 are progressively decreasing their contribution to the group. Hence, this analysis may help the group coordinator to evaluate researchers' work and quantitatively justify his/hers actions. Moreover, the centroid position is a fundamental methodology to analyze the specific members that engage in the collective work.

(4) Conclusions

Our work aims to present a methodology to analyze and evaluate research groups and cooperation between researchers. Thus, a MatLab script was developed which benefits from networking concepts described in the literature, as well as others introduced specifically in this work, aiming to analyze the performance of inter and intraresearchers' performance. Using the herein proposed methodology in the multi-disciplinary research group RoboCorp, it was possible to observe an increasing level of cooperation among researchers over the years, thus indicating an increase of the individual and collective intervention towards the success of the group. It was also possible to observe the existence of dependencies between researchers, as well as hierarchical classifications within the group. The results led to a further understanding of the dynamical cooperation between researchers, emphasizing the importance of this methodology for the practical evaluation of research groups, thus enhancing potential useful relations between researchers. As future work, it is considered important to include other weighting factors, as the impact factor of articles and the number of authors associated to each publication. It is also relevant to analyze the researchers who contribute the most to the success of the group. To that end, new metrics need to be explored to further evaluate the usefulness of the work developed by researchers and their real contribution and relevance to the collective performance.

Acknowledgments

This work was supported by a PhD scholarship (SFRH/BD /73382/2010) granted to the first author by the Portuguese Foundation for Science and Technology (FCT); Institute of Systems and Robotics (ISR); RoboCorp and Instituto de Telecomunicações (IT-Covilhã); and also the work was under regular funding by FCT.

References

- Albert R; Jeong H; and Barabasi AL (2000) Error and Attack Tolerance of Complex Networks. *Nature*, **406** (6794): 378–382. Available at http://www.nature.com/nature/ journal/v406/n6794/full/
- Arunachalam S; Srinivasan R; and Raman V (1994) International Collaboration in Science: Participation by the Asian Giants. Scientometrics, 30 (1): 7-22.
- **Boorman SA** (1975) A Combinatorial Optimization Model for Transmission of Job Information through Contact Networks. *The Bell Journal of Economics*, **6** (1): 216-249.
- Brueckner JK; and Spiller PT (1991) Competition and Merges in Airline Networks. *International Journal of Industrial Organization*, **9** (3): 323-342.
- Cardillo A; Scellato S; and Latora V (2006) A Topological Analysis of Scientific Coauthorship Networks. *Physica A*, **372** (2): 333-339.
- Cardoso J; Mendling J; Neumann G; and Reijers HA (2006) A Discourse on Complexity of Process Models. *Lecture Notes in Computer Science*, **4103**: 117-128.

- Ernesto Estrada (1995) Edge Adjacency Relationships in Molecular Graphs Containing Heteroatoms: A New Topological Index Related to Molecular Volume. *Journal of Chemical Information and Computer Science*, 35 (4): 701-707.
- Fiduccia CM; and Mattheyses R M (1982) A Linear-Time Heuristic for Improving Network Partitions. In: Proceeding of Design Automation 19th Conference, IEEE, 14-16 June 1982, General Electronic Research and Development Center, Schenectady, Las Vegas NY, USA, pp175-181.

Available at http://www.ieeexplore.ieee.org

Fortunato S; and Castellano C (2009) Community Structure in Graphs. In; Santo Fortunsto etal, (eds), Complex Networks, Proceedings of ComplexNet 2009, International Workshop on Complex Networks, 26-27 May 2009, Catania, Italy, pp1-42.

Available at: http://www.sites.google.com/ site/santofortunato/publications2

- Fortunato S (2010) Community Detection in Graphs. *Physics Reports*, **486**: (3-5) 75-174. Available at http://www.adsabs.harvard.edu/ abs/2010OhR...486...75F
- Freeman L (1978) Centrality in Social Networks: Conceptual Clarification. *Social Networks*, 1: 215–239.

Available at: http://www.moreno.ss.nci. edu/27pdf

- Harzing AW (2007) Publish or Perish. *Research in International Management Products and Services for Academics.* Associate Dean Research, University of Melbourne, Australia. Available at: http://www.harzing.com/pop.htm.
- Hendricks K; Piccione M; and Tan G (1995) The Economics of Hubs: The Case of Monopoly. *Review of Economic Studies*, 62 (1): 83-99. Available at: http://www.links.jstor.org/ sici?sici=0034-6527%28199501%
- Hespanha JP (2004) An Efficient MATLAB Algorithm for Graph Partitioning. Technical Report, University of California, Oct. 2004, California, USA.

Available at: http://www.ece.ucsb.edu

Horvath S (2011) Weighted Network Analysis: Applications in Genomics and Systems Biology. Springer, 1st ed. Springer, New York, USA, pp1-446.

- Katz M L; and Shapiro C (1994) Systems Competition and Network Effects. *Journal of Economic Perspectives*, 8 (2): 93-115.
- Keren M; and Levhari D (1983). The Internal Organization of the Firm and the Shape of Average Costs. *The Bell Journal of Economics*, 14 (2): 474-486.
- Kleinberg J (1999) Authoritative Sources in a Hyperlinked Environment. *Journal of the ACM*, **46** (5): 604-632.
- Lim C; Bohacek S; Hespanha J; and Obraczka K (2005) Hierarchical Max-Flow Routing. In Proceedings of the IEEE GLOBECOM. Available at: http://www.eecis.udel.edu/
- Narin F; Stevens K; and Whitlow ES (1991) Scientific Co-operation in Europe and the Citation of Multinationally Authored Papers. *Scientometrics*, **21** (3): 313-323.
- Newman MEJ (2004a). Co-authorship Network and Patterns of Scientific Collaboration. Proceedings of the National Academy of Sciences of the United States of America, 101 (Supl. 1): 5200-5205.full

Available at: http://www.pnas.org/

- Newman MEJ (2001) The Structure of Scientific Collaboration Networks. *Proceedings of the National Academy of Sciences of the United States of America*, **98** (2): 404-409. Available at: http://www.pnas.org/
- Newman MEJ; Watts DJ; and Strogatz SH (2002) Random Graph Models of Social Networks. Proceedings of the National Academy of Sciences of the United States of America, 99 (Supl. 1): 2566-2572.

Available at: http://www.pnas.org/contents/99/ sup1/2566-2572.full

Newman, MEJ (2004b) Who is the best Connected Scientist? A Study of Scientific Co-authorship Networks. *Lectures Notes in Physics*, **650**: 337-370.

Available at: http://www.springer.com/ chaptr/10.1007%02F970-3-540-44485-5_16#page-1/

- **Okubo Y; Miquel JF; Frigoletto L;** and **Dore JC** (1992) Structure of International Collaboration in Science: Typology of Countries through Multivariate Techniques using a Link Indicator. *Scientometrics*, **25** (2): 321-351.
- Otte E; and Rousseau R (2002) Social Network Analysis: a Powerful Strategy, also for the Information Sciences. *Journal of Information Science*, **28** (6): 441-453.

- Rodriguez MA; and Pepe A (2008) On the Relationship between the Structural and Socioacademic Communities of a Coauthorship Network. *Journal of Informetrics*, 2 (3): 195-201.
- Ravasz E; Somera AL; Mongru DA; Oltvai ZN; and Barabasi AL (2002) Hierarchical Organization of Modularity in Metabolic Networks. *Science*, **297** (5586): 1551–1555. Available at: http://www.sciencemag.org/ contents/297/5586/1551.abstract
- Watts DJ; and Strogatz SH (1998) Collective Dynamics of 'Small-World' Networks. *Nature*, **393** (6684: 440-442. Available at http://www.nature.com/nature/
- Watts DJ (2002) A simple Model of Global Cascades on Random Networks. *Proceedings* of the National Academy of Sciences of the United States of America, **99** (9), 5766–5771. Available at: http://www.pnas.org/ contents/99/9/5766–5771.full
- Wu M (2009) wgPlot-Weighted Graph Plot. MatLab Central File Exchange. *MATLAB Central*,