# Evaluating the impact of interdisciplinary research: a multilayer network approach

Elisa Omodei, Manlio De Domenico and Alex Arenas

Department of Mathematics and Computer Science, Rovira i Virgili University

(Dated: January 25, 2016)

Nowadays, scientific challenges usually require approaches that cross traditional boundaries between academic disciplines, driving many researchers towards interdisciplinarity. Despite its obvious importance, there is a lack of studies on how to quantify the influence of interdisciplinarity on the research impact, posing uncertainty in a proper evaluation for hiring and funding purposes. Here we propose a method based on the analysis of bipartite interconnected multilayer networks of citations and disciplines, to assess scholars, institutions and countries interdisciplinary importance. Using data about physics publications and US patents, we show that our method allows to reveal, using a quantitative approach, that being more interdisciplinary causes – in the Granger sense – benefits in scientific productivity and impact. The proposed method could be used by funding agencies, universities and scientific policy decision makers for hiring and funding purposes, and to complement existing methods to rank universities and countries.

## I. INTRODUCTION

Interdisciplinary research has recently gained a central role in the advancement of science, leading to important achievements [1]. For instance, the 2014 Nobel Prize in Chemistry was awarded to two physicists and a physical chemist, for "a physical technique, developed with help from chemistry, that helps illuminate problems in biology" [2].

Even though several definitions and metrics for interdisciplinarity have been proposed [3–6], citation impact metrics accounting for this aspect of scientific research have not been defined yet.

On the other hand, funding agencies have created specific calls for interdisciplinary projects, like the Interdisciplinary Programs funded by the National Science Foundation [7]. The European Research Council explicitly encourages applications from scientists having published in multidisciplinary journals [8], and the evaluation criteria for the Marie Curie fellowships also include the interdisciplinary aspects of the research [9]. Consequently, there is significant need to evaluate projects and scholars by considering interdisciplinarity too. The difficulties in evaluating interdisciplinary research constitute a pressing controversy that leads many young scholars to remain on more traditional tracks, because the risks associated to undertaking an interdisciplinary career path seem too high [10].

This work addresses the issue of quantifying interdisciplinarity by proposing a method to rank scientific publications (such as papers and patents) and their producers (scholars, inventors, institutions, companies and countries) according to their scientific impact and its breadth over different scientific disciplines. The method is based on the detection of the most central elements of a complex bipartite interconnected multilayer network representing scientific producers and scientific citations within and across different fields. The citation network is composed of multiple layers, each representing a scientific discipline. Accounting for this diversity – instead of neglecting the information it provides by building an aggregated representation of the network – allows to unveil the cross-disciplinary versatility of scientific publications and of their producers, and therefore to obtain a quantitative measure of their interdisciplinary scientific impact.

Since the seminal work of de Solla Price [11] and Garfield [12], scientists have put a great effort into trying to understand the patterns of citation distributions [13–15] and the non-trivial dynamics of scientific recognition [16–24]. This fundamental body of work has the ultimate goal of setting the basis for the definition of more accurate and fairer scientific impact metrics used for evaluation purposes.

In the last decades, several indices have been presented. They are based on the idea that we can quantify the impact of a scientific publication by counting the number of citations it has received over the years. A widely adopted indicator to evaluate scholars' scientific impact is the h-index [25] (a scholar has an index $h$ if $h$ of her/his publications have received at least $h$ citations each), and its numerous variants [26–28].

More recently, a different approach has been proposed. Building networks that reconstruct the chains of scientific citations allows for a global understanding of the intrigued patterns of citations between publications – or between producers. This representation allows to unveil the difference between, for instance, a publication that has received 10 citations coming from highly cited publications, and a publication that has received 10 citations too, but from low-cited papers. The two have the same number of citations but the former has clearly had a higher impact. To rank publications, journals or scholars according to their importance in the respective citation network, researchers have proposed diffusion algorithms that simulate the spreading of scientific credits on the network [29–31]. In practice, this is the same idea at the basis of the PageRank, i.e the algorithm that Google uses to rank the pages of the World Wide Web [32].
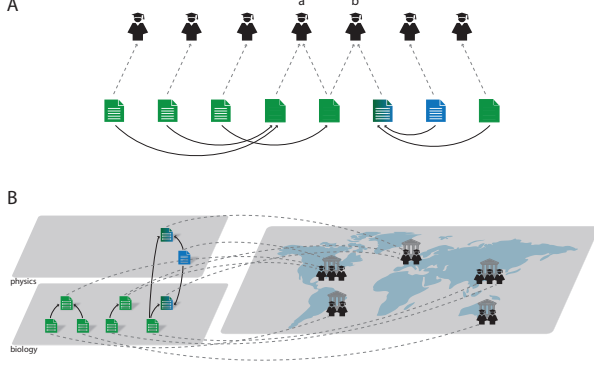
Figure 1: **Bipartite interconnected multilayer network.** Panel A shows a simple example of bipartite citation network made of 8 papers and 7 scholars. The 8 papers belong to two disciplines – biology and physics. Green icons represent biology papers, blue physics, and the bicolour icon represents a paper that belongs to both biology and physics. Continuous arrows represent citation edges, whereas dotted arrows connect papers to its authors. Panel B shows the multilayer representation of the network. Consider, for example, authors $a$ and $b$. If we discard the information about the scientific fields and consider the aggregated network shown in panel A, then the two authors' centrality would be the same, because they authored the same number of papers, having an identical structure of incoming citations. However, the multilayer framework takes into account that one of $b$'s papers pertains to both physics and biology, and, moreover, had an impact in both fields (one citation comes from a physics paper and the other from a biology one). Therefore $b$ has a higher versatility than $a$.

## II.  METHODOLOGY

In this work, we propose to rank scientific publications and their producers employing the PageRank defined on a bipartite interconnected multilayer structure that accounts for citations within and across different disciplines. This is equivalent to ranking nodes according to their *versatility* [33] on an interconnected multilayer network [34].

To account for interdisciplinarity, we define a bipartite interconnected multilayer network representing citations between publications (papers or patents) and relations between publications and their manufacturers (scholars, inventors, research institutions, companies or countries).

Given $N$ nodes and $L$ layers, the rank$-4$ multilayer adjacency tensor $A_{\beta\tilde{\delta}}^{\alpha\tilde{\gamma}}$ is defined in the following way. Let $C_\beta^\alpha(\tilde{h}, \tilde{k}) = \sum_{i,j=1}^{N} w_{i,j}(\tilde{h}, \tilde{k}) E_\beta^\alpha(ij)$ be the rank$-2$ adjacency tensor encoding information about the relationship between layer $\tilde{h}$ and $\tilde{k}$, where $w_{ij}(\tilde{h}, \tilde{k})$ indicates the intensity of the relationship between node $n_i$ in layer $\tilde{h}$ and node $n_j$ in layer $\tilde{k}$, and $E_\beta^\alpha(ij)$ indicates the rank$-2$ tensor that represents the canonical basis in the space

$\mathbb{R}^{N \times N}$ (note that when $\tilde{h} = \tilde{k}$, $C_\beta^\alpha(\tilde{h}, \tilde{h})$ represents the intra-layer adjacency tensor), then

$$A_{\beta\tilde{\delta}}^{\alpha\tilde{\gamma}} = \sum_{\tilde{h}, \tilde{k}=1}^{L} C_\beta^\alpha(\tilde{h}, \tilde{k}) E_{\tilde{\delta}}^{\tilde{\gamma}}(\tilde{h}, \tilde{k}) \qquad (1)$$

where $E_{\tilde{\delta}}^{\tilde{\gamma}}(\tilde{h}, \tilde{k})$ indicates the rank$-2$ tensor that represents the canonical basis in the space $\mathbb{R}^{L \times L}$. This is the general formulation of an adjacency tensor representing a multilayer network.

To build our network we consider $N = N_P + N_A$ nodes (where $N_P$ is the number of publications, and $N_A$ the number of manufacturers of the chosen type that produced the $N_P$ papers. Therefore, given the ordered set of nodes $\{n_1, ..., n_N\}$, the first $N_P$ elements $\{n_1, ..., n_{N_P}\}$ represent publication, and the other $N_A$ elements $\{n_{N_P+1}, ..., n_N\}$ represent manufacturers. Moreover, we consider $L = L' + 1$ layers, where $L'$ is the number of scientific disciplines that the publications belong to. The 4 components of the rank$-2$ adjacency tensor $C_\beta^\alpha(\tilde{h}, \tilde{k})$ are defined as follows. $C_\beta^\alpha(\tilde{l}_x, \tilde{l}_x)$ and $C_\beta^\alpha(\tilde{l}_x, \tilde{l}_y)$, with $x, y \in [1, L']$, encode information about publication citations. Each layer represents a discipline or a subfield, therefore $w_{ij}(\tilde{l}_x, \tilde{l}_x) = \frac{1}{N_L(i)N_L(j)}$ if both publications $i$ and $j$ belong to discipline $x$, and publication $i$ cites publication $j$. $N_L(i)$ ($N_L(j)$) is the number of disciplines that publication $i$ ($j$) belongs to. This normalisation is performed so that every citation carries one unit of value overall. Interdisciplinary citations are instead encoded by $C_\beta^\alpha(\tilde{l}_x, \tilde{l}_y)$; $w_{ij}(\tilde{l}_x, \tilde{l}_y) = \frac{1}{N_L(i)N_L(j)}$ if publications $i$ belongs to discipline $x$ and $j$ to discipline $y$, and publication $i$ cites publication $j$. Let $\tilde{l}_A$ denote the remaining layer, then the tensors $C_\beta^\alpha(\tilde{l}_x, \tilde{l}_A)$, with $x \in [1, L']$, encode information about the relation between publications and their manufacturers, *i.e.* if the chosen type of manufacturer is scholars, then $w_{ij}(\tilde{l}_x, \tilde{l}_A) = \frac{1}{N_L(i)}$ if author $j$ is one of the authors of publication $i$. If we consider research institutions, we connect each publication to the institutions to which its authors are affiliated; if we consider countries, the connections are to the countries in which these institutions are based. Finally, we define $C_\beta^\alpha(\tilde{l}_A, \tilde{l}_x)$ and $C_\beta^\alpha(\tilde{l}_A, \tilde{l}_A)$ to be zero tensors. $C_\beta^\alpha(\tilde{l}_A, \tilde{l}_x)$ tensors are null because we do not want the relations between publications and manufacturers to be symmetric, to avoid unrealistic paths to take place when computing the nodes centrality. $C_\beta^\alpha(\tilde{l}_A, \tilde{l}_A)$ is null because all the information is already encoded in the other tensors: we do not need to explicitly add citation edges between authors.

In this framework, for the citation layers, each node is active on a given layer if and only if the publication it represents belongs to the corresponding field. For example, a monodisciplinary publication is active only on one layer, whereas an interdisciplinary publication, pertaining to both physics and biology, is active on two layers. As a consequence, a publication whose impact is restricted to only one discipline has intra-layer incoming edges only,

whereas a publication that has influenced the work of researchers in more than one field has inter-layer incoming edges too, which represent the bridges between the different fields involved. Therefore this framework allows for a natural representation of the interdisciplinarity degree of a publication. Being our goal to rank publication producers too, we introduce in the network a second type of nodes, which, according to the specific need, represent scholars, inventors, research institutions, or countries. These nodes are active on a dedicated layer, and each publication has directed outgoing inter-layer edges pointing to each of its producers. Previous works on ranking producers are based on one-mode projections of the bipartite network of publications and producers, whereas in this work we prefer to take advantage of the complete bipartite structure in order to avoid any information loss, as further detailed in Appendix A.

On the proposed network, the ranking is obtained through a process of diffusion of scientific credits from paper to paper through citation edges within and across disciplines. Producers are the sinks of this diffusion process, being represented by nodes with no outgoing edges, and incoming edges originated from the papers they have produced. A schematic representation of the proposed network is shown in Figure 1.

Having defined the multilayer citation network, we propose to rank its nodes according to their PageRank versatility, which is given by the steady-state solution of the equation

$$p_{\beta\tilde{\delta}}(t+1) = R^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}} p_{\alpha\tilde{\gamma}}(t) \tag{2}$$

where $p_{\alpha\tilde{\gamma}}(t)$ is the time-dependent tensor that gives the probability to find a random walker at a particular node $\alpha$ in a particular layer $\tilde{\gamma}$, and

$$R^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}} = \left[ r T^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}} + \frac{1-r}{NL} u^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}} \right], \tag{3}$$

$N$ being the number of nodes in the network, $L$ the number of layers, and $r$ the teleportation rate. $T^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}}$ denotes the rank$-4$ tensor of transition probabilities for jumping between pairs of nodes and switching between pairs of layers, and $u^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}}$ is the rank$-4$ tensor with all components equal to 1. Let $\Omega_{\alpha\tilde{\gamma}}$ be the eigentensor of the transition tensor $R^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}}$, denoting the steady-state probability to find a random walker in node $\alpha$ and layer $\tilde{\gamma}$. The tensor $\Omega_{\alpha\tilde{\gamma}}$ provides the PageRank of each node ($\alpha$) in each layer ($\tilde{\gamma}$): it is crucial to remark here that this is not equivalent to calculate the PageRank in each layer separately, because our formulation accounts for the whole interconnected structure to solve the eigenvalue problem. To obtain the multilayer PageRank of each node, regardless of the layer, we project the values obtained from its replicas in different layers, obtaining the multilayer PageRank vector

$$\omega_\alpha = \Omega_{\alpha\tilde{\gamma}} u^{\tilde{\gamma}} \tag{4}$$

where $u^{\tilde{\gamma}}$ is the vector with all components equal to 1. It has been shown [33] that this operation provides the same results that would be obtained by calculating PageRank by means of simulated random walkers that explore the multilayer structure according to transition rules encoded in $R^{\alpha\tilde{\gamma}}_{\beta\tilde{\delta}}$.

## III. DATA

To illustrate the proposed ranking method, we test it on two case studies: the APS and the US patents datasets.

The first is a collection of papers published in the journals of the American Physical Society (Physical Review Letters, Physical Review and Review of Modern Physics) between 1985 and 2009 [?  ]. We restricted the analysis only to papers with at most ten authors, to avoid biases due to the papers of experimental high-energy physics in which all the project collaborators are listed as co-authors. To disambiguate author's name, we used a simple technique introduced in previous studies [31]. Metadata in the dataset provide information about the topic of the papers through the specification of the assigned "Physics and Astronomy Classification Scheme" (PACS) code, developed by the American Institute of Physics (AIP) and used in Physical Review since 1975 to identify fields and sub-fields of physics [35]. We exploited this information to build a heterogeneous interconnected 10-layer network in which each layer represents a subfield of physics, as defined by the PACS systems: *General*; *The Physics of Elementary Particles and Fields*; *Nuclear Physics, Atomic and Molecular Physics*; *Electromagnetism, Optics, Acoustics, Heat Transfer, Classical Mechanics, and Fluid Dynamics*; *Physics of Gases, Plasmas, and Electric Discharges*; *Condensed Matter: Structural, Mechanical and Thermal Properties*; *Condensed Matter: Electronic Structure, Electrical, Magnetic, and Optical Properties*; *Interdisciplinary Physics and Related Areas of Science and Technology*; *Geophysics, Astronomy, and Astrophysics*. From the paper meta-data we also extracted the authors affiliation information, which allowed us to associate to each paper a list of (one or more) institutions and countries. The final dataset consists of 319816 papers, 204809 authors, 626 institutions and 54 countries. Arguably, the APS cover only Physics, but note that physics is a vast field that spans from biological physics to astrophysics and although it may fall short of a full interdisciplinary analysis it is clear that this is a powerful indicator of multi-topic analysis that serves to proof the usefulness of the method.

The second dataset contains the U.S. patents granted between January 1963 and December 1999, and all citations made to these patents between 1975 and 1999 [36]. To define the layers, we used the 6 categories proposed in previous studies [37]: *Chemical (excluding Drugs)*; *Computers and Communications*; *Drugs and Medical*; *Electrical and Electronics*; *Mechanical*; *Others*. Each patent

is assigned to one main class defined by the United States Patent and Trademark Office (USPTO), and to any number of subsidiary classes. Each class belongs to one of the listed categories, therefore each patent is associated with one or more layer according to its classes. However, the dataset only contains the information about the main class, therefore we complemented it by extracting the information about the other classes from the USPTO Patent Grant Full Text [38] The final dataset contains 1574882 patents, 1142499 inventors, 138833 assignees (i.e. corporations for the most part), and 127 countries.

## IV. RESULTS

Figure 2 shows the evolution of the interdisciplinary ranking of the world top physics departments, and of the world top companies, over time. This visualization allows to observe, for instance, the raise of the University of Texas at Austin during the 1990s, after the establishment, in 1985, of the Center for Nonlinear Dynamics, funded and directed by the Boltzmann Medal laureate Harry Swinney [39].

Using the same data, we show that the proposed method is able to capture two fundamental aspects of interdisciplinary research: intrinsic multidisciplinarity (i.e. publishing papers or patents pertaining to different areas) on the one hand, and effective interdisciplinarity, i.e. being credited by different scientific areas, on the other.

We define the *topical interdisciplinarity* $TI(a)$ of an author $a$ (who could be a scholar or an inventor) as the average number of different scientific areas her/his publications pertain to, i.e.:

$$TI(a) = \frac{1}{n(a)} \sum_{i=1}^{n(a)} d(p_i) \qquad (5)$$

where $n(a)$ is the number of publications authored by $a$, and $d(p_i)$ is the number of fields that publication $p_i$ belongs to.

Moreover, for each publication $p$ we define an entropy metrics based on the distribution of its incoming citations across the different fields represented by the different layers:

$$H(p) = \sum f_i \log \frac{1}{f_i} \qquad (6)$$

where the sum is over the different fields (layers) and $f_i$ is the proportion of edges incident in $p$ that come from layer $i$. Therefore if a publication is only cited by other publications belonging to its own field $H(p) = 0$, whereas a publication that has received citations from different fields has $H(p) > 0$, and the higher the number of fields it has had impact on, the higher its entropy. For each author $a$ we then compute her/his *citation interdisciplinar-*
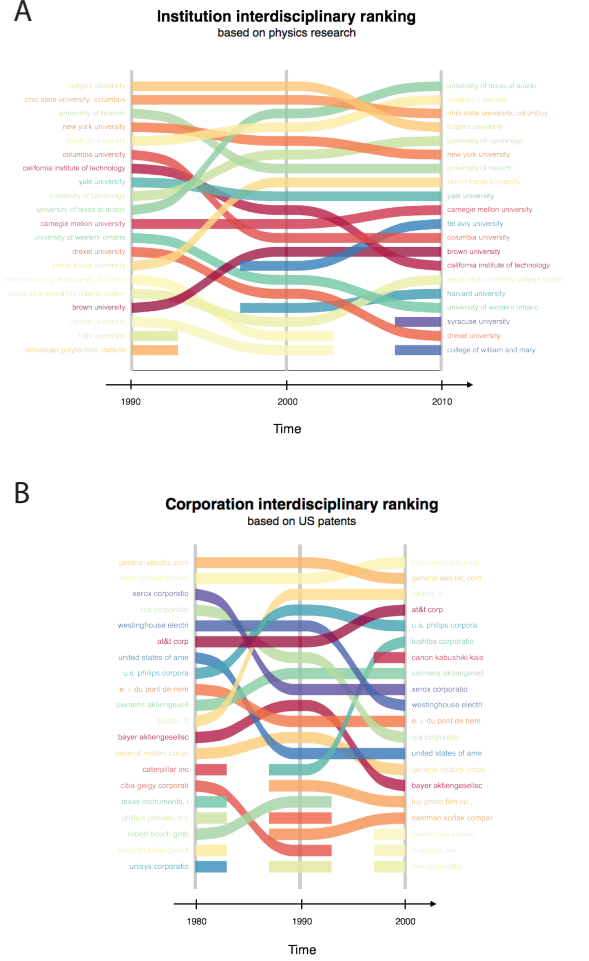


Figure 2: **Interdisciplinary ranking evolution.** Panel A: Visualization of the time evolution of the interdisciplinary impact ranking of the top 20 physics departments, computed using the APS dataset. Panel B: Time evolution of the interdisciplinary impact ranking of the top 20 world companies, computed using the US patent data. Broken lines represent institutions or companies that do not belong to the top 20 in the previous or the following time stamp.

*ity* $CI(a)$ as the average entropy of her/his publications:

$$CI(a) = \frac{1}{n(a)} \sum_{i=1}^{n(a)} H(p_i). \qquad (7)$$

We find a strong positive correlation between the gain in rank that scholars and inventors obtain when evaluated using the proposed method – instead of a method based on a flat representation of the citation network –, and their topical interdisciplinarity (Figure 3, panels (a) and (b)). Moreover, we find that the rank gain is also positively correlated with the disciplinary diversity of scholars' and inventors' incoming citations (Figure 3, panels (c) and (d)).

Recent studies have shown correlations between pa-

APS journals authors     US patents inventors

r = 0.75 (CL = 99.9%)     r = 0.76 (CL = 99.9%)
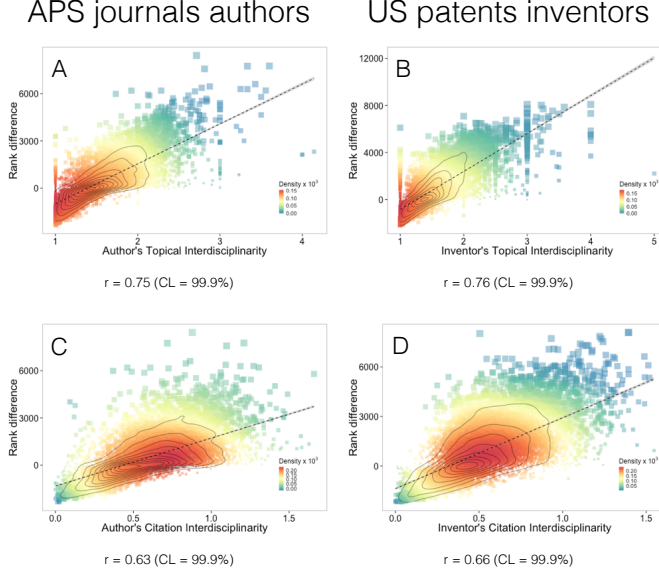
r = 0.63 (CL = 99.9%)     r = 0.66 (CL = 99.9%)

Figure 3: **Correlations.** Heat-maps representing the correlation between the gain in rank that scholars and inventors obtain when evaluated using the proposed method – instead of a method based on a flat representation of the citation network –, and two measures of their interdisciplinarity level. The x-axis represents, in panel (a) and (b), scholars' and inventors' topical interdisciplinarity, defined as the average number of different scientific areas their publications pertain to, and, in panel (c) and (d), their diversity in terms of disciplines of the scholars' and inventors' incoming citations (citation interdisciplinarity). Correlations are calculated using Pearson's $r$ coefficient, and setting the statistical significance at 0.1%. Solid lines represent density gradient contours, and dashed lines represent linear regression models estimated via maximum-likelihood.

pers' and producers' interdisciplinarity and their citation impact [40, 41]. Here, having defined a quantitative method to rank scientific research actors according to their interdisciplinary impact, we carry out a longitudinal study based on Granger-causality of the implications of interdisciplinarity on scientific research. We find that the trend of an institution' or company' interdisciplinary impact over time Granger-causes [42] its trend in productivity and impact. To obtain this result we check, for every institution in the APS dataset and every company in the US patent dataset, over a 20 years time span, if the evolution of the interdisciplinary impact Granger-causes the evolution of the corresponding h- and r- index [43]. The latter is defined as the square root of the sum of the citations to the publications in the h-core, i.e. the set of $h$ publications having received at least $h$ citations each. It has been shown that these two indices complement each other in evaluating the production and impact of scientific publication producers [27].

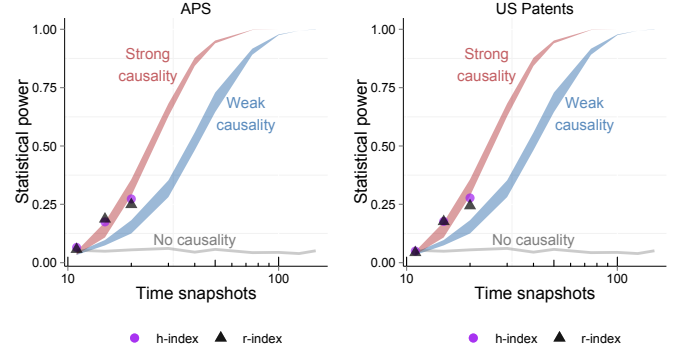Given the reduced size of the time series, only 20



Figure 4: **Granger-causality.** Plot representing the statistical power of Granger-causality tests in function of the number of data points available, using real data (circles and triangles), causal time series toy models (red and blue lines), and null models of non-correlated time series (grey line). Line width represents the range of variability of the model across different lags. We can observe how the statistical power of the real data series is highly consistent with the one obtained using strong causal time series with the same number of data points. This indicates that interdisciplinary impact has causal repercussions – in the Granger sense – on scientific productivity and impact.

temporal observations, we first need reference values to asses whether we can accept the hypothesis that the observed time series are Granger-causal. We generate 1000 strongly causal autoregressive models models with memory of the same size of the data, and quantify the corresponding statistical power $\pi$, i.e. the fraction of times we correctly accept the alternative hypothesis when it is true (being the null hypothesis that the observed time series are *not* Granger-causal). We find a value of $\pi \sim 0.3$, with a 95% confidence level obtained through bootstrapping [44]. We perform the same test on 1000 non-causal series and find $\pi \sim 0.05$, almost an order of magnitude lower. Therefore a value of $\pi$ that lies in the interval between these two reference values indicates Granger-causality (see Appendix B for further details). Indeed, applying the same analysis on the two empirical datasets, we find that the interdisciplinary impact time series Granger-causes both the h- and r- index time series with a statistical power $\pi = 0.28$ and $\pi = 0.25$, respectively (Figure 4). Even if we can not quantify the exact amount of Granger-causality, because of the lack of time series spanning a longer period, we can conclude that our study provides quantitative evidence, with 95% confidence level, that there is a causal relationship between interdisciplinarity, productivity and impact.

## V. DISCUSSION

In this paper, we propose a methodology to assess the citation impact of scientific publications and their

producers that intrinsically accounts for their interdisciplinarity. This aspect was not included in previous citation impact indicators.

Even though numerous metrics have been proposed to assess citation impact, several issues have been raised. These include the accounting of self-citations [45], the choice of the appropriate citation time window [46], field normalisation [47], and author credit allocation [48]. Despite the vast literature on the subject, consensus is still lacking on how to solve these issues. Here, we propose a method whose objective is to account for interdisciplinarity, and we do not enter those debates. However, the bipartite interconnected multilayer networks of citations and disciplines that we introduce can be adapted to take into account specific needs. For example, edges connecting papers to their authors could be weighted differently to take into account non-homogeneous allocation of credit, or a specific time window could be chosen a priori to select the papers constituting the network.

Going beyond the presented assessment of the benefits produced by interdisciplinarity, the method proposed in this work could constitute a tool for funding agencies and academic hiring decision makers to quantify the impact of interdisciplinary research and its producers, for a faster advancement of excellent science.

### Acknowledgements

### Appendix A: Comparison with other approaches

The idea of ranking scholars simulating a diffusion process was already introduced in [31] – and previously in [29] to rank scientific publications – but the approach proposed in this work considers a different kind of network – a bipartite interconnected multilayer network. In this section we motivate the choice of taking into account the complete bipartite structure instead of its one-mode projection.

For the sake of simplicity, we will consider a bilayer version of the network. Our focus here is not in fact the multilayer aspect of the network capturing interdisciplinarity, but rather its bipartition. Therefore, let us consider $N = N_P + N_A$ nodes and 2 layers $\{l_1, l_2\}$. The 4 components of the rank$-2$ adjacency tensor $C_\beta^\alpha(\tilde{h}, \tilde{k})$ are now defined as follows. $C_\beta^\alpha(\tilde{l}_1, \tilde{l}_1)$ encodes information about citing relations between papers, i.e. $w_{ij}(\tilde{l}_1, \tilde{l}_1) = 1$ if paper $i$ cites paper $j$. $C_\beta^\alpha(\tilde{l}_1, \tilde{l}_2)$ encodes information about paper authorship, i.e. $w_{ij}(\tilde{l}_1, \tilde{l}_2) = 1$ if author $j$ is one of the authors of paper $i$. Finally, we define $C_\beta^\alpha(\tilde{l}_2, \tilde{l}_1)$ and $C_\beta^\alpha(\tilde{l}_2, \tilde{l}_2)$ to be zero tensors, consistently with the representation introduced in S.1. For the sake of simplicity, since in the rest of the section we will be dealing only with rank-2 tensors, we will make use of the simpler classical

matrix notation instead of the tensorial one. Therefore we will denote $C_\beta^\alpha(\tilde{l}_1, \tilde{l}_1)$ as $C$ and $C_\beta^\alpha(\tilde{l}_1, \tilde{l}_2)$ as $A$.

In the author citation network proposed in [31], each node represents an author, and $w_{ij} \neq 0$ if there exist at least one publication $\alpha$, of which $i$ is an author, that cites a publication $\beta$ of which $j$ is an author. Each such publication gives a contribution $\frac{1}{nm}$ to $w_{ij}$ (where $n$ is the number of authors of publication $\alpha$, and $m$ is the number of authors of $\beta$) so that the total contribution of each citation is equal to 1.

Let us consider an adjacency matrix $\mathcal{C}$ of size $N_P \times N_P$, encoding the citation links between papers. $C$ can be built from $\mathcal{C}$ by means of multiplication with a rectangular matrix $\mathcal{I}$ of size $(N_P + N_A) \times N_P$ such that $(\mathcal{I})_{ii} = 1$ for $i = 1, ..., N_P$, and all the other elements are equal to 0. Then

$$C = \mathcal{I}\mathcal{C}\mathcal{I}^T. \tag{A1}$$

Using $\mathcal{I}$ we can also build $A$ from the projection matrix $\mathcal{P}$, of size $N_P \times N_A$, where $w_{ij} = 1$ if $j$ is one of the authors of paper $i$:

$$A = \mathcal{I}\mathcal{P}\mathcal{I}^T. \tag{A2}$$

Let us define $\tilde{\mathcal{P}}$ as the normalised version of $\mathcal{P}$, i.e. $(\tilde{\mathcal{P}})_{ij} = \frac{(\mathcal{P})_{ij}}{\sum_{k=1}^{N_A}(\mathcal{P})_{ik}} = \frac{1}{m_i}$ (where $m_i$ is the number of authors of paper $i$), then the $N_A \times N_A$ adjacency matrix representing the network of citations between authors can be obtained performing two successive matrix multiplications:

$$\mathcal{A} = \tilde{\mathcal{P}}^T \mathcal{C} \tilde{\mathcal{P}}. \tag{A3}$$

*Proof:*

$$(\tilde{\mathcal{P}}^T\mathcal{C})_{ik} = \sum_{h=1}^{N_P}(\mathcal{P})_{hi}(\mathcal{C})_{hk} = \sum_{\substack{h=1 \\ (\mathcal{P})_{hi} \neq 0, (\mathcal{C})_{hk}=1}}^{N_P} \frac{1}{m_h} \tag{A4}$$

This means that $(\tilde{\mathcal{P}}^T\mathcal{C})_{ik}$ is a sum over the papers authored by $i$ that cite paper $k$, where each paper $h$ gives a contribution of 1 over the number of authors. Then:

$$(\mathcal{A})_{ij} = \sum_{k=1}^{N_P}(\tilde{\mathcal{P}}^T\mathcal{C})_{ik}(\mathcal{P})_{kj} = \sum_{\substack{k=1 \\ (\mathcal{P})_{kj} \neq 0}}^{N_P}\Big(\sum_{\substack{h=1 \\ (\mathcal{P})_{hi} \neq 0, (\mathcal{C})_{hk}=1}}^{N_P} \frac{1}{m_h}\Big)\frac{1}{m_k} \tag{A5}$$

Each element $(\mathcal{A})_{ij}$ is therefore a sum over all the pairs of papers $(h, k)$ such that $i$ is an author of $h$ and $j$ of $k$, and each element of the sum gives a contribution equal to $\frac{1}{m_h m_k}$, as indeed defined in [31]. $\square$

We have demonstrated that the adjacency matrix of the author citation network is obtained performing two operations of matrix multiplication involving $\mathcal{C}$ and $\mathcal{P}$. Matrix multiplications consists in multiplications and summations of the matrix elements, which inevitably lead to information loss. The supra-adjacency matrix of the network proposed in this paper is instead the sum of the two expansions of $\mathcal{C}$ and $\mathcal{P}$, i.e. $C$ and $A$, respectively. This guarantees that no information is loss, and this why we chose to consider the whole bipartite structure. Figure 5 shows an example in which the information loss characteristic of the author citation network leads to a less fair ranking of authors compared to the ranking based on the network introduced in this paper. Using our approach (top figure), the most central author is $B$, who is the author of both the most central paper and of one of the second most central papers. However, in the author citation network

framework, the most central author is $A$ (to understand why, we recall that the PageRank centrality is based not only on the number of incoming edges, but on the importance of the nodes from which these edges originate). This is due to the fact that in the author citation network all the information about an author's incoming citations from different papers is aggregated, and therefore $A$ benefits from the importance of $B$ without any distinction between the importance coming from papers that actually cite $A$, and that coming from $B$'s other papers. In this case $B$'s most central (and cited) paper is not the one citing $A$'s paper, but this information is lost in the author citation network. As a consequence, the resulting ranking does not always reflect the real importance of the different authors.
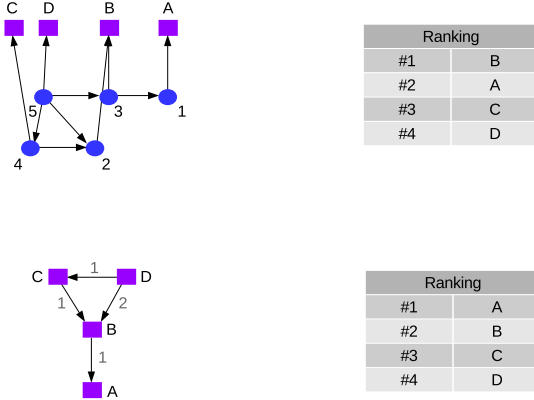


Figure 5: An example in which the PageRank centralities computed on the author citation network and on the bipartite network lead to different rankings of the authors.

An alternative approach is to use the PageRank method to get the centrality of papers in $\mathcal{C}$, and then compute the author centrality as a properly normalised sum of the centralities of the papers she/he has authored. In matrix terms, the author PageRank centrality vector $\omega_A$ can be obtained by simply applying a linear transformation to the paper PageRank centrality vector $\omega_P$:

$$\omega_A = \tilde{\mathcal{P}}^T \omega_P. \tag{A6}$$

However, this solution involves another kind of aggregation which can lead to misleading results too. An example is shown in Figure 6. Using the sum approach, author $D$ becomes more central than author $A$, because she/he authored two papers, even though they are the two most marginal papers in the network (note that the only citation to paper 4 is a self-citation). On the contrary, $A$ is the author of a very central paper, and in fact our approach correctly classifies her/him as more central than $D$. The issue with this alternative approach is that the PageRank is a diffusion process, which is not a linear dynamics. Therefore summing over the centralities of different nodes is also an aggregation process through which some information on the system is lost.
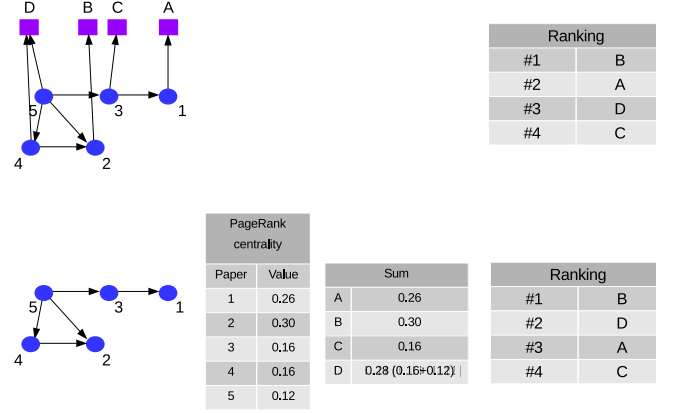


Figure 6: An example in which the PageRank centralities computed as the sum over the papers and on the bipartite network lead to different rankings of the authors.

## Appendix B: Testing Granger-causality using bootstrapping

To test whether a given time series $x(t) = [x(1), \ldots, x(T)]$ Granger-causes a second time series $y(t) = [y(1), \ldots, y(T)]$, where $T$ is the number of available measurements, we use the following procedure. Using the likelihood-ratio test statistics, we test the null hypothesis of non-Granger-causality against the alternative hypothesis of Granger-causality, for three different values of the maximum lag (1,2, and 3). In fact, we do not have any *a priori* knowledge of the lag the Granger-causality could depend on. We denote with $\hat{D}^\star$ the maximum value obtained for the likelihood-ratio. Successively, we build $n$ surrogate time series which are not Granger-causal by randomly shuffling $x(t)$ and $y(t)$, performing the tests again and obtaining the values $\{\hat{D}_1, \hat{D}_2, ..., \hat{D}_n\}$. To give our results statistical relevance, we use a two-sided nonparametric Wilcoxon-Mann-Whitney rank test that, in our case – one empirical observation against several synthetic values – can be summarized as follows. Given the value of the test statistics obtained from the data, we compare it against the values obtained from null models and if $\hat{D}^\star$ is larger or smaller than $\hat{D}_i$, for *all* $i = 1, 2, ..., n$, than we reject the null hypothesis. This condition corresponds to verify iff

$$\hat{D}^\star < \min_i \{\hat{D}_i\} \quad \text{OR} \quad \hat{D}^\star > \max_i \{\hat{D}_i\}. \tag{B1}$$

To obtain a confidence level of 95% (i.e., type I error $\alpha = 0.05$), we need to generate $n = \frac{2}{\alpha} = 40$ surrogates. We follow this procedure for all the available time series and then compute the Granger-causality statistical power as the fraction of series for which the alternative hypothesis is accepted.

To evaluate the statistical power that we could expect for time series with only 20 data points, we generate different sets of synthetic models where we know, *a priori*, which hypothesis is true (the null or the alternative one). First, we generate causal time series using autoregressive models with memory, while varying the number of temporal observations from 11 to 150. In particular, we generate 1000 series for each one of

the following models:

$$y[t] = 0.75x[t-l] + \epsilon \text{ (strongly causal models)} \quad \text{(B2)}$$

$$y[t] = 0.5y[t-1] + 0.5x[t-l] + \epsilon \text{ (causal models)} \quad \text{(B3)}$$

with $l$ ranging in the interval $[1,3]$, and $\epsilon$ being white noise. For these models we know that the alternative hypothesis is true.

Then we generate 1000 time series which are not Granger-causal:

$$y[t] = 0.5y[t-1] + \epsilon. \quad \text{(B4)}$$

In all cases the equation governing the $x$ time series is:

$$x[t] = 0.5x[t-1] + \epsilon. \quad \text{(B5)}$$

For these models we know that the null hypothesis is true.

We perform the procedure described above for each model separately. In the cases when the null hypothesis is true, for a test with 95% confidence level we expect to wrongly reject the null hypothesis (i.e., to accept the alternative hypothesis) only 5% of times, exactly what we observe as shown in Figure 3 of the main text. In the cases when the alternative hypothesis is true, the fraction of times we correctly accept it depends, as expected, on the size of the time series, where the availability of a larger number of temporal observations leads to a higher rate of correct decisions. For time series with 20 points, the maximum achievable power in the classification of causal time series is $\approx 0.3$, in good agreement with results obtained from the empirical datasets (Figure 4 of the main text).

[1] Nature, Nature **525**, 289 (2015).
[2] ChemistryWorld, Chemistry World (2014).
[3] A. L. Porter, A. S. Cohen, J. D. Roessner, and M. Perreault, Scientometrics **72**, 117 (2007).
[4] L. Leydesdorff, Journal of the American Society for Information Science and Technology **58**, 1303 (2007).
[5] C. S. Wagner, J. D. Roessner, K. Bobb, J. T. Klein, K. W. Boyack, J. Keyton, I. Rafols, and K. Börner, Journal of Informetrics **5**, 14 (2011).
[6] P. Jensen and K. Lutkouskaya, Scientometrics **98**, 619 (2014).
[7] *https://www.nsf.gov/od/iia/additional_resources/interdisciplinary_research/support.jsp* (2015).
[8] *"applicants should also be able to demonstrate a promising track-record of early achievements appropriate to their research field and career stage, including significant publications (as main author) in major international peer-reviewed multidisciplinary scientific journals, or in the leading international peer-reviewed journals of their respective field." http://erc.europa.eu/funding-and-grants/funding-schemes/starting-grants* (2015).
[9] *Annex 2 http://ec.europa.eu/research/participants/portal/doc/call/h2020/h2020-msca-if-2015/1645199-guide_for_applicants_if_2015_en.pdf* (2015).
[10] D. Rhoten and A. Parker, Science(Washington) **306**, 2046 (2004).
[11] D. J. de Solla Price, Science **149**, 510 (1965).
[12] E. Garfield and R. K. Merton, *Citation indexing: Its theory and application in science, technology, and humanities*, vol. 8 (Wiley New York, 1979).
[13] S. Redner, The European Physical Journal B-Condensed Matter and Complex Systems **4**, 131 (1998).
[14] D. A. King, Nature **430**, 311 (2004).
[15] F. Radicchi, S. Fortunato, and C. Castellano, Proceedings of the National Academy of Sciences **105**, 17268 (2008).
[16] R. Guimera, B. Uzzi, J. Spiro, and L. A. N. Amaral, Science **308**, 697 (2005).
[17] M. Newman, EPL (Europhysics Letters) **86**, 68001 (2009).
[18] Y. Eom, S. Fortunato, and M. Perc, PLoS ONE **6**, e24926 (2011).
[19] D. Wang, C. Song, and A.-L. Barabási, Science **342**, 127 (2013).
[20] O. Penner, R. K. Pan, A. M. Petersen, K. Kaski, and S. Fortunato, Scientific reports **3** (2013).
[21] Q. Zhang, N. Perra, B. Gonçalves, F. Ciulla, and A. Vespignani, Scientific reports **3** (2013).
[22] B. Uzzi, S. Mukherjee, M. Stringer, and B. Jones, Science **342**, 468 (2013).
[23] P. Deville, D. Wang, R. Sinatra, C. Song, V. D. Blondel, and A.-L. Barabási, Scientific reports **4** (2014).
[24] Q. Ke, E. Ferrara, F. Radicchi, and A. Flammini, Proceedings of the National Academy of Sciences p. 201424329 (2015).
[25] J. E. Hirsch, Proceedings of the National academy of Sciences of the United States of America **102**, 16569 (2005).
[26] L. Egghe, Scientometrics **69**, 131 (2006).
[27] L. Bornmann, R. Mutz, and H.-D. Daniel, Journal of the American Society for Information Science and Technology **59**, 830 (2008).
[28] J. Kaur, F. Radicchi, and F. Menczer, Journal of Informetrics **7**, 924 (2013).
[29] D. Walker, H. Xie, K.-K. Yan, and S. Maslov, Journal of Statistical Mechanics: Theory and Experiment **2007**, P06010 (2007).
[30] C. Bergstrom, College & Research Libraries News **68**, 314 (2007).
[31] F. Radicchi, S. Fortunato, B. Markines, and A. Vespignani, Physical Review E **80**, 056103 (2009).
[32] S. Brin and L. Page, Computer networks and ISDN systems **30**, 107 (1998).
[33] M. De Domenico, A. Solé-Ribalta, E. Omodei, S. Gómez, and A. Arenas, Nature communications, in press (2015).
[34] M. De Domenico, A. Solé-Ribalta, E. Cozzo, M. Kivelä, Y. Moreno, M. A. Porter, S. Gómez, and A. Arenas, Physical Review X **3**, 041022 (2013).
[35] *http://journals.aps.org/PACS* (2015).
[36] *http://www.nber.org/patents/* (2015).
[37] B. H. Hall, A. B. Jaffe, and M. Trajtenberg, Tech. Rep., National Bureau of Economic Research (2001).
[38] *http://www.google.com/googlebooks/uspto-patents-grants-text.html* (2015).
[39] *https://web2.ph.utexas.edu/utphysicshistory/UTexas_Physics_History/Center_for_Nonlinear_Dynamics.html* (2015).
[40] V. Larivière, S. Haustein, and K. Börner, PloS one **10** (2015).

[41] L. Pan and S. Katrenko, *A Review of the UK's Interdisciplinary Research using a Citation-based Approach* (Elsevier, 2015).

[42] C. W. Granger, Econometrica: Journal of the Econometric Society pp. 424–438 (1969).

[43] B. Jin, L. Liang, R. Rousseau, and L. Egghe, Chinese science bulletin **52**, 855 (2007).

[44] B. Efron and R. J. Tibshirani, *An introduction to the bootstrap* (CRC press, 1994).

[45] W. Glänzel, K. Debackere, B. Thijs, and A. Schubert, Scientometrics **67**, 263 (2006).

[46] J. Wang, Scientometrics **94**, 851 (2013).

[47] Y. Li, F. Radicchi, C. Castellano, and J. Ruiz-Castillo, Journal of Informetrics **7**, 746 (2013).

[48] M. Gauffriau, P. Larsen, I. Maye, A. Roulin-Perriard, and M. von Ins, Scientometrics **77**, 147 (2008).