



# Ground truth? Concept-based communities versus the external classification of physics manuscripts

Vasyl Palchykov<sup>1,2\*</sup>, Valerio Gemmetto<sup>1</sup>, Alexey Boyarsky<sup>1</sup> and Diego Garlaschelli<sup>1</sup>

\*Correspondence:

palchykov@lorentz.leidenuniv.nl

<sup>1</sup>Lorentz Institute for Theoretical Physics, Leiden University, Niels Bohrweg, 2, Leiden, 2333CA, The Netherlands

<sup>2</sup>Institute for Condensed Matter Physics, Svientsitskii str. 1, Lviv, 79011, Ukraine

## Abstract

Community detection techniques are widely used to infer hidden structures within interconnected systems. Despite demonstrating high accuracy on benchmarks, they reproduce the external classification for many real-world systems with a significant level of discrepancy. A widely accepted reason behind such outcome is the unavoidable loss of non-topological information (such as node attributes) encountered when the original complex system is converted to a network. In this article we systematically show that the observed discrepancies may also be caused by a different reason: the external classification itself. For this end we use scientific publication data which (i) exhibit a well defined modular structure and (ii) hold an expert-made classification of research articles. Having represented the articles and the extracted scientific concepts both as a bipartite network and as its unipartite projection, we applied modularity optimization to uncover the inner thematic structure. The resulting clusters are shown to partly reflect the author-made classification, although some significant discrepancies are observed. A detailed analysis of these discrepancies shows that they may carry essential information about the system, mainly related to the use of similar techniques and methods across different (sub)disciplines, that is otherwise omitted when only the external classification is considered.

**Keywords:** science of science; community detection; bipartite networks

## 1 Introduction

A conflict between two members of a relatively small university organization that happened more than 40 years ago [1] has attracted a lot of attention in the scientific community so far [2]. A confrontation during the conflict resulted in a fission of the organization, known as Zachary's karate club, into two smaller groups, gathered around the president and the instructor of the club, respectively. Predicting the sizes and compositions of the resulting factions, given the structure of the social interaction network before the split, attracted a lot of attention. This puzzle, supplemented by the known outcome, makes this system among the best studied benchmarks to test community detection algorithms [3]. Having verified a high level performance on the aforementioned system and on other benchmarks [4], community detection algorithms have then been massively applied to uncover tightly connected modules within large real-world systems. This allowed scientists

to identify, for instance, Flemish- and French-speaking communities in Belgium using mobile phone communication networks [5], detect functional regions in the human or animal brain from neural connectivity [6], observe the emergence of scientific disciplines [7] and investigate the evolution of science using citation patterns and article metadata [8–10].

A bird's eye view on the identified clusters in real-world systems certifies their meaningfulness. However, an in-depth quantitative validation of the community structure requires its comparison with an external classification of the nodes, which is accessible only for a limited number of large systems. Examples include crowd-sourced tag assignments for software packages [11], product categories for Amazon copurchasing networks [12], declared group membership for various online social networks [13, 14] and publication venues for co-authorship networks in the computer science literature [13]. Surprisingly, significant discrepancies have been identified between the extracted grouping of nodes and their external classification for these systems [11, 15]. This message remains robust independently of the system under investigation and the technique used to uncover its community structure, and calls for a detailed inspection of such discrepancies in order to understand the reasons behind them.

One of the possible reasons concerns the strong simplification that occurs during the projection of the original complex system into a network. This projection may omit some crucial information that cannot be encoded into the structural connection pattern [11]. The missing information may correspond to age or gender of individuals in social networks [16, 17] or geographical position of the nodes within spatially embedded systems [18]. Following this direction, several algorithms [19, 20] have been developed in order to handle specific nodes attributes, beside the usual connectivity patterns. Such approaches have been shown to identify groups of nodes that more closely reproduce the external classification in real-world systems [20] than the techniques that rely on the connectivity patterns only.

In this article we argue that, independently of the aforementioned issue, the supposedly poor performance of community detection algorithms may be caused by the external classification itself and its misinterpretation. For instance, a system may possess several alternative classification schemes, such as thematic and methodological groupings in a system of scientific publications or in academic co-authorship networks [21]. In such situation, the discrepancies between the community detection results and a single accessible classification (e.g. based on thematic similarity) may carry, instead, meaningful information (e.g. about methodological similarity), therefore providing an added value to the system understanding.

In this article we explore this idea by performing a detailed analysis of a scientific publication record system. This system may be simplified to structural network representation, where the nodes correspond to scientific articles, and the links represent the relationship between them. There are various possibilities to map these relationships: direct citation [22], co-citation and bibliographic coupling [23] or content related similarities [24, 25]. Here we focus on the latter, considering scientific terms or concepts that appear within the articles. Performing community detection on the corresponding network, we compare the results with an expert made classification of these articles, considering both similarities and discrepancies between the two different partitions. Then we investigate the main reasons causing the most notable deviations.

This article is organized as follows. In the Data section we present the dataset used; in Methods we introduce the methodology used to build the networks, extract the partitions and compare them with the external classification. Finally, in Results and Conclusions we present our findings and discuss them.

## 2 Data

We investigate a collection of scientific manuscripts submitted to e-print repository arXiv [26] during the years 2013 and 2014. During the submission process, the authors were requested to classify the manuscript according to the arXiv classification scheme by assigning at least one category to it. In our analysis we are focussed only on the articles that have been assigned to a single category, restricting ourself to the field of physics. Moreover, the collections of manuscripts submitted during the years 2013 and 2014 are considered separately, eliminating the possible issues related to the temporal evolution of research disciplines. The resulting datasets consist of 36,386 articles submitted during 2013 and 41,848 articles submitted during 2014, and will be referred below (together with the extracted contents) as the arxivPhys2013 and arxivPhys2014 datasets, respectively. The numbers of articles belonging to each category are shown in Table 1.

Each article is represented by a set of scientific concepts that characterize its content, i.e. specific words or combinations of them. The concepts have been identified within the full text by the ScienceWISE.info platform (SW). SW is a web service connected to the main online repositories such as arXiv, whose peculiarity is a bottom-up approach in the management of scientific concepts [27]. The initially created scientific ontology was followed by a continuous editing by the users, for instance by adding new concepts, definitions and relationships. This crowd-sourced procedure leads to the most comprehensive vocabulary of scientific concepts in the domain of physics. Such vocabulary takes care of synonyms that refer to the same concepts and it includes physics concepts explicitly labeled as generic like mass or energy, or more specific ones like community detection. Both are the results of crowd-sourcing by the registered expert-users.

**Table 1** Distribution of articles among categories

Category	$n_{2013}^s$	$n_{2013}^m$	$n_{2014}^s$	$n_{2014}^m$
nucl-th	648	1,628	766	1,210
nucl-ex	315	924	324	736
hep-ph	2,625	3,935	3,116	2,885
hep-ex	602	1,726	706	1,225
hep-lat	356	695	419	417
hep-th	1,787	3,717	2,316	2,960
gr-qc	1,118	2,782	1,527	2,204
astro-ph	10,984	3,023	11,445	2,437
physics	4,452	6,479	5,711	4,880
cond-mat	10,549	4,609	11,397	3,538
nlin	392	327	522	905
quant-ph	2,558	3,240	3,187	2,471
math-ph	0	3,789	412	2,668

The number of manuscript submitted during the year  $y$  that have been assigned to a given category only ( $n_y^s$ ) or to the category and at least one another ( $n_y^m$ ). List of categories: theoretical and experimental nuclear physics (nucl-th and nucl-ex, respectively), four branches of high energy physics (hep-ph: phenomenology, hep-ex: experiment, hep-lat: lattice and hep-th: theory), general relativity and quantum cosmology (gr-qc), astrophysics (astro-ph), physics (physics), condensed matter physics (cond-mat), nonlinear science (nlin), quantum physics (quant-ph) and mathematical physics (math-ph).

The number  $k$  of concepts significantly vary among the manuscripts, reaching up to  $k_{\max} \sim 400$  for review articles. The average number of identified concepts  $\langle k \rangle$  per article, together with some other characteristics of the datasets `arxivPhys2013` and `arxivPhys2014`, are shown in Table 2. The datasets supporting the conclusions of this article are included within Additional file 1.

### 3 Methods

The dataset may be represented as a network, whose nodes correspond to articles. Two nodes  $i$  and  $j$  are connected by a link if the corresponding articles share at least a single common concept. The resulting networks are extremely dense, covering almost 90% of all possible network connections; this number may be reduced to 50% if the generic concepts are ignored (see Table 2). Below, to save the computational resources, we will ignore the generic concepts in our analysis. The weight of the link between two nodes is designed to reflect the level of content similarity between two articles, i.e. the overlap between the respective lists of concepts. Different concepts, however, may contribute differently to the similarity among two articles. Indeed, sharing a widely used concept should affect the similarity between two articles differently than sharing a specific one, suggesting that specific concepts should have a higher impact on the similarity. Each concept  $c$  in the dataset is therefore weighted according to its occurrence, which may be accounted for by the so-called `idf(c)` factor [28]:

$$\text{idf}(c) = \log \frac{N}{N(c)}. \tag{1}$$

Here  $N$  is the total number of articles and  $N(c)$  is the number of articles that contain concept  $c$ . As mentioned above, among the  $V$  concepts identified by `SW`, we will consider only the specific ones, discarding the  $V_{\text{gen}}$  generic concepts. The content of each article can be therefore expressed by means of a  $(V - V_{\text{gen}})$ -dimensional concept vector  $\vec{v}_i$ . The element  $v_{ic}$  of the concept vector of the article  $i$  has non-zero value equal to `idf(c)` only if the concept  $c$  appears within the article  $i$  and equals zero otherwise.

The similarity between the contents of two articles  $i$  and  $j$ , and the link weight  $w_{ij}$  between the corresponding nodes, may then be estimated by the cosine similarity between the two concept vectors  $\vec{v}_i$  and  $\vec{v}_j$  as follows:

$$w_{ij} = \frac{\vec{v}_i \cdot \vec{v}_j}{|\vec{v}_i| |\vec{v}_j|}. \tag{2}$$

The resulting network will be referred below as the `idf` representation of the data.

**Table 2 Basic characteristics of the datasets**

	$N$	$V$	$V_{\text{gen}}$	$\langle k \rangle$	$L_{\text{idf}}^{\text{in}}$	$L_{\text{idf}}$	$L_{\text{bp}}^{\text{in}}$	$L_{\text{bp}}$
<code>arxivPhys2013</code>	36,386	12,200	347	37	$5.9 \times 10^8$	$3.3 \times 10^8$	$2.1 \times 10^6$	$1.3 \times 10^6$
<code>arxivPhys2014</code>	41,848	12,728	344	38	$7.8 \times 10^8$	$4.5 \times 10^8$	$2.5 \times 10^6$	$1.6 \times 10^6$

Total number of articles ( $N$ ), total number of identified concepts ( $V$ ) and the number of generic ones ( $V_{\text{gen}}$ ) among them;  $\langle k \rangle$  gives the average number of non-generic concepts within arbitrary chosen article. The number of links in a unipartite network provided that the generic concepts are included ( $L_{\text{idf}}^{\text{in}}$ ) or excluded ( $L_{\text{idf}}$ ) is two orders of magnitude larger than the corresponding number of links in bipartite networks ( $L_{\text{bp}}^{\text{in}}$  and  $L_{\text{bp}}$ , respectively). This results in significant differences in computational resources needed to perform community detection analysis.

Alternatively to *idf* representation, the dataset may be mapped into a bipartite network. Such network consists of the nodes of two types that correspond to manuscripts and scientific concepts, respectively. The unweighted links in the simplest case reflect the appearance of a concept within the article. This network will be referred below to as a *bp* representation of the data, and the usage of the two alternative representation will serve the robustness of our results. The number of links ( $L_{idf}$ ,  $L_{bp}$ ) of these networks are shown in Table 2. As one may see, the number of links in *bp* representation is about two orders of magnitude smaller than the number of links in the corresponding *idf* representation. This has significant consequences on the run-time and memory used to analyse the networks.

Indeed, the run-time  $t$  of the employed algorithm [5] scales about linearly with the number of links  $L$  of the considered network. Since empirically in the bipartite representation  $L_{bp} \sim O(N)$  while in the unipartite case  $L_{idf} \sim O(N^2)$ , this reflects in much different computational resources required to perform the community detection. Moreover, here we point out that the bipartite representation is the most natural and suitable characterization of the dataset, since the null model behind such representation of the data is definitely more correct. In fact, the bipartite null model is consistent with the constraints on both the types of node (number of papers per concept and concepts per article). This feature is instead lost when the system is projected into a unipartite network, since the previous constraints are not matched any more. Furthermore, the bipartite representation and null model already take into account the presence of more frequent concepts, sparing us the use of any *idf* factor. In this context, we therefore propose the use of the bipartite representation as a possible alternative to the more widespread *idf* (or *tf-idf*) unipartite representation.

In order to find a unipartite network partition, we will maximize a modularity function [29]. To deal with bipartite networks, we adopt a co-clustering approach [30] and Barber's generalization of modularity [31].

In both cases, we assume that each article may belong to a single cluster only, hence exploiting the notion of non-overlapping communities. Furthermore, the co-clustering approach makes stronger restrictions on a bipartite partition, compared to a unipartite one. Indeed, the resulting clusters of a bipartite partition consist of both articles and related concepts, and we assume that each concept belongs to a single cluster as well. Such restriction may be relaxed, for instance by using alternative ways to generalize modularity for bipartite network [32] or by employing stochastic block model techniques [33]. However, we will consider co-clustering of bipartite networks since it allows us to straightforwardly employ the same greedy optimization algorithm [5] for the networks of both types.

The restriction towards a single algorithm is also caused by the result [11] that (i) the selected algorithm is among the ones that perform best on real-world networks and (ii) the major influence on the accuracy is related to the dataset itself rather than the algorithm. Due to the stochastic origin of this algorithm, it has been applied 100 times for unipartite networks and 1,000 times for bipartite ones (due to significantly different number of links and, therefore, the required computational resources). Among the detected partitions, for each network we will select the single partition that corresponds to the highest value of modularity; this partition will be referred below as the optimal partition for each network.

#### 4 Results

A partition of a bipartite network consists of clusters that contain both articles and scientific terms (concepts), while clusters of a unipartite network partition consist of articles only. To compare both unipartite and bipartite partitions with the external article classification, we will be focussed only on the articles that fall into each cluster. Thus, by referring below to a cluster of bipartite partition we mean the set of articles that belong to the specified cluster. In this perspective, the external classification of the articles is represented by the arXiv standard split into different subject classes or categories (astro-ph, cond-mat, etc.).

Then, given two partitions  $P$  and  $Q$  of the same network (for instance a detected network partition and the arXiv classification), an initial comparison between them has been performed using an information-based symmetrically normalized mutual information:

$$I_N(P, Q) = \frac{2I(P, Q)}{H(P) + H(Q)}. \tag{3}$$

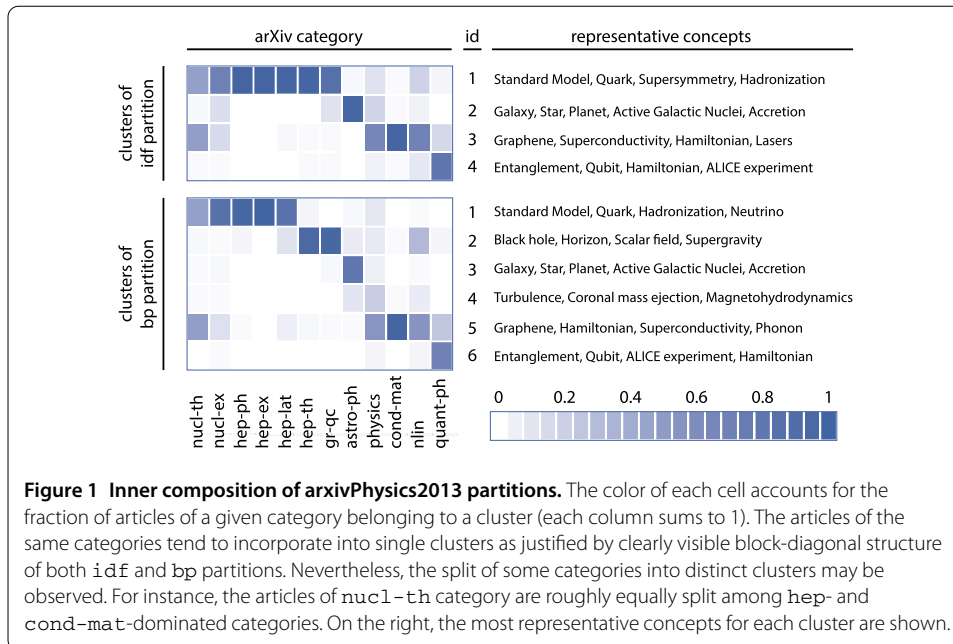
Here  $I(P, Q)$  is the mutual information [34] between two partitions  $P$  and  $Q$ , and  $H(P)$  is the entropy of partition  $P$ . The normalized mutual information  $I_N(P, Q)$  may vary between 0 and 1. A value of 0 indicates that the two partitions have no information in common, while a value of 1 corresponds to identical partitions. In Table 3 we show the level of similarity between the resulting partitions and the arXiv classification ones. The reported values of normalized mutual information indicate the existence of some common information between automatically identified clusters of articles (both in the bipartite and unipartite cases) and the author based classification. However, the values being quite far from the possible maximum of 1 reflect evidence for some discrepancies between the partitions. Below we perform a detailed analysis of these discrepancies and show the results for the arxivPhys2013 dataset. Similar findings can be observed in the arxivPhys2014 case and they are shown in Additional file 2.

The first difference is observed in the numbers of detected clusters and of arXiv subject classes: while the number of categories in the arXiv classification scheme is 12,<sup>a</sup> the number of clusters in our partitions is only equal to 4 in the idf and to 6 in the bp network representations, respectively.<sup>b</sup> Indeed, the articles of some different arXiv categories tend to belong to a single cluster. This may be clearly observed in Figure 1 that shows the fraction of articles of each arXiv category belonging to each cluster in the resulting partitions. This merger is especially visible for different high energy physics (hep) categories (hep-ph, hep-ex, hep-lat and hep-th): in the idf partition, almost 99% of all these articles fell into a single cluster, independently of the sub-field. This result, despite deviating from the arXiv classification scheme, is reasonable since we observe a union of almost all papers about high energy physics, no matter if they deal with experimental or theoretical issues.

**Table 3 Similarity between network partitions and external classification**

	idf	bp
arxivPhys2013	0.600 ± 0.025	0.563 ± 0.026
arxivPhys2014	0.553 ± 0.002	0.536 ± 0.023

Average value of the normalized mutual information  $I_N$  (3) between a partition of each network representation and arXiv classification of the articles and the corresponding standard deviations. Both bp and idf partitions demonstrate similar value of closeness to arXiv classification.

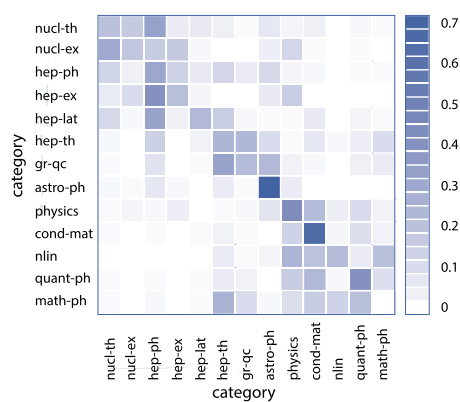


Instead, in the *bp* partition the articles of the four *hep* categories are almost entirely distributed among two clusters, focussed on experimental and theoretical issues, respectively. The first of them joins 95% of all articles that belong to experimental categories (*hep-ph*, *hep-ex* or *hep-lat*), while the second one contains 94% of all theoretical (*hep-th*) articles. Thus, the presence of more clusters within the bipartite network partition allows us to identify methodologically different clusters of articles within the *hep* categories, in particular dividing theoretical papers from experimental ones.

Even though the split of *hep* articles into two groups may be simply explained by the different approaches used to study the phenomena, a further result can be observed from Figure 1: in the bipartite network partition, *hep-th* articles tend to form a single cluster with the articles that belong to general relativity and quantum cosmology (category *gr-qc*) rather than with the other high energy physics articles, thus appearing to be more similar to *gr-qc* papers rather than to the other *hep* ones. Intuitively, indeed, we know that both *hep-th* and *gr-qc* both focus mostly on general relativity, while the other *hep* categories focus on particle physics.<sup>c</sup>

Such relatedness between the articles of the two theoretical physics categories (*hep-th* and *gr-qc*) may be verified independently by a category co-occurrence analysis. To show this, we will use the complementary part of the investigated dataset. This set consists of all articles that have been submitted to *arXiv* during the same 2013 year, but for which the authors have assigned at least two different categories. Thus, no article of this set overlaps with the clustered *arxivPhys2013* collection. Irrespective of the details of the decision-making process through which authors assign multiple categories, this multiplicity reflects the author’s decision that the scope of the article can not be properly covered by a single category of a given classification scheme. Whilst several categories may cover the scope of a single research article, the co-occurrence of the same two categories in a significant fraction of articles may reflect some hidden relationships between them. The corresponding empirical co-occurrence matrix is shown in Figure 2 and indicates the fraction of articles of a given category that have been co-submitted to the other categories. The diagonal

**Figure 2 Co-occurrence matrix of arXiv categories during year 2013.** Built on the complementary dataset to `arxivPhys2013`, this matrix reflects the relationships between arXiv categories and allows to justify the meaningfulness of some remarkable discrepancies, like the merger of `hep-th` and `gr-qc` articles. Each non-diagonal element reflects the fraction of articles in which two specified categories have co-occurred. The diagonal cells represent the fractions of articles that have been assigned a single category, i.e. they concern the articles of the `arxivPhys2013` dataset. A normalization procedure has been performed such that each row of the matrix sums to 1. Thus, the aforementioned fractions correspond to the fractions of manuscripts that have been labeled with a given category.



elements of this matrix indicate the fraction of articles of each category that have been assigned a single category by the author(s), i.e. the articles of the `arxivPhys2013` dataset. A normalization procedure has been performed such that each row of the matrix sums to 1.

Figure 2 confirms that the `hep-th` subject class is indeed more related to the `gr-qc` class than to the other `hep` categories: `hep-th` co-occurred with `gr-qc` in 1,721 articles, and with all other `hep` categories in only 1,286 articles, even though the number of the corresponding `hep` papers (`hep-ph`, `hep-ex`, `hep-lat`) exceeds the number of `gr-qc` ones threefold. This high level of relatedness between `hep-th` and `gr-qc` categories justifies the merging of the articles of these categories into a single cluster and indicates the meaningful deviation from the arXiv classification scheme. It is worth to mention that in the `idf` partition, where all `hep` category articles tend to belong to a single cluster, the same cluster is supplemented by 87% of all `gr-qc` articles, in agreement with the result observed above. Moreover such a tendency is not restricted to the dataset for the selected year: it has also been observed for the `arxivPhys2014` one.

The same approach explains the presence of a significant fraction of `physics`, non-linear (`nlin`) and quantum physics (`quant-ph`) articles in `cond-mat` clusters. It also allows us to understand a possible reason why nuclear physics articles (both theory and experiment) occur significantly within `hep` clusters. However, it cannot explain the presence of roughly one half of `nucl-th` articles in the condensed matter cluster (cluster no. 3 in `idf` and no. 5 in `bp` partitions) in both network representations. The latter deviation from the article classification, which is not explained by category co-occurrence, does not exclude that similarities between these topics exist but are considered not strong enough by the authors to label the articles with both subject classes. To uncover the possible essence of these similarities, we examine the top representative concepts that characterize the `nucl-th` articles that belong to the two different clusters, see Table 4. In both cases, the top representative concepts contain the ones that characterize the object of investigation within theoretical nuclear physics, such as `Isotope`, `Isospin` or `Nuclear matter`. However, one may clearly identify method-related concepts, such as `Hartree-Fock`, `Hamiltonian` and `Mean field`, among the top representative concepts of articles in the `cond-mat` cluster. These concepts clearly characterize methods that are widely used in condensed matter physics research, and that have not been identi-



**Table 4** Top representative concepts of two groups of articles categorized as `nuc1-th`

%	Concept (cluster no. 1)	%	Concept (cluster no. 3)
43	Hadronization	55	Isotope
39	Isospin	53	Hamiltonian
37	Pion	39	Hartree-Fock
33	Degree of freedom	36	Quadrupole
32	Heavy ion collision	34	Isospin
31	Quark	31	Nuclear matter
29	Chirality	30	Degree of freedom
29	Hamiltonian	28	Mean field
29	Nuclear matter	26	Harmonic oscillator
26	Coupling constant	25	Spin orbit

The left side of the table represents the group of articles that fell into `hep` dominated cluster (no. 1) in `idf` partition. The right side - the other group: the `nuc1-th` articles that fell into `cond-mat` dominated cluster (no. 3). For each group, the numbers next to the concepts give the percentage of articles in which the concept has been identified. The table allows us to make a suggestion that the two groups of articles significantly differ by the methods used to investigate nuclear matter.

fied among top concepts in any other cluster. This result emphasizes the ability of scientific concepts found within research articles to highlight not only topics focussed on the same objects, but also methodologically similar research directions.

## 5 Conclusions

The differences between the outcomes of community detection algorithms and possible external classifications may have various reasons. The most notable of them concern a possible failure of the considered algorithm or the unavoidable loss of data about real complex systems determined by their representation as networks. To deal with the first issue, algorithms are heavily tested on benchmarks, while the second issue is still under investigation [20]. In this article, we emphasize a third possible reason behind such discrepancies, i.e. the fact that the external classification itself may possess its own limitations. For this reason we performed a detailed investigation of a scientific publication records, which (i) may be naturally represented as a network and (ii) owns an external author-made classification of articles. While, indeed, some discrepancies are caused by the lack of data (for instance in the case of the articles for which no concept has been identified), we argue that the most remarkable of them may reflect real commonalities across different subject classes. Academic publications are traditionally categorized and classified<sup>d</sup> according to objects or phenomena under investigation. The same phenomena, however, may be explored using various approaches, experimental observation and theoretical modeling being among them. On the other hand, the phenomena that belong to different research topics may be investigated using the same methods, composing the core of the interdisciplinary research. Thus, a more comprehensive classification of research articles may be represented by a two layer categorization scheme, where one layer reflects phenomena or objects while the other one stands for the methods of investigation. Usually, these two layers are not taken equally into account. The expert made classification may include rather a strong bias towards the object layer. The reasons involve the classification scheme itself and the limited knowledge about all other research disciplines that employ the same methods. Instead, automatic concept-based categorization should have no direct preference for any of the layers: the extracted concepts correspond both to phenomena and methods, and the algorithm has no information about the possible division of the concepts. Thus, the observed discrepancies may reflect the dominance of the methodological layer over the other one, which corresponds to phenomena or objects.

Similar results have been previously observed within the collaboration network of scientists at Santa Fe Institute [21], where, besides the expected grouping around common topics, some methodologically driven clusters have been observed.

This shows that the failure in reproducing an external classification may indicate a genuinely more complicated organization within the system, in addition to the lack of data or algorithmic mistakes. Besides developing sophisticated algorithms to deal with real systems, we should therefore keep in mind that some observed discrepancies may go beyond the standard classification and carry important information about the system under study. We believe that similar results may be observed in other systems. Indeed, the ground truth necessarily follows from a given classification criterion; however, the considered data may contain more than that single type of information (perhaps in conflict one with each other). In general, therefore, it may happen that what we consider as the ground truth is just one of the possible reference points, rather than some absolute truth. Understanding the information employed to define the so-called ground truth is therefore crucial in order to perform a proper comparison between external classification and automatically retrieved communities.

## Additional material

**Additional file 1: Data sets.** The data file contains used metadata together with the lists of extracted concept identifiers for each manuscript under investigation. (zip)

**Additional file 2: Community detection results for articles submitted during year 2014.** The figure consists of the inner composition of `idf` and `bp` partitions of `arXivPhys2014` dataset and the corresponding category co-occurrence matrix. These results to a large extent reproduce the results obtained for the year 2013, thus verifying the conclusions made. (pdf)

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

All authors discussed and designed the experiments as well as contributed to the writing of the paper. VP and VG implemented and conducted the experiments. All authors read and approved the final manuscript.

## Acknowledgements

The authors thank A Cardillo, A Martini, P de Los Rios, O Ruchaiskiy and D Larremore for useful discussions, and A Magalich for preparation of the data. This work was supported by SNSF project No. 147609 Crowdsourced conceptualization of complex scientific knowledge and discovery of discoveries and by the EU project MULTIPLEX (contract 317532).

## Endnotes

- <sup>a</sup> In fact, there are 13 physics categories in `arXiv` classification scheme, but there is no single article in `arXivPhys2013` dataset that belong to `math-ph` category only.
- <sup>b</sup> By performing a detailed comparison we ignore all single-node clusters, which contain the articles for which no concepts has been identified.
- <sup>c</sup> Indeed, it is very likely that nowadays the `hep` categories would be split in multiple subcategories (namely `hep-th`, `hep-lat`, etc.). However, here we point out that our study (in particular in the bipartite case) shows that `hep-th` looks actually more similar to `gr-qc` than to the other `hep-` classes. This therefore seems to strengthen the apparently counterintuitive choice of dividing the high energy articles in different primary classes.
- <sup>d</sup> Document classification and categorization are different processes: classification refers to the assignment one or more predefined categories to a document, while categorization refers to the process of dividing the set of documents into priory unknown groups whose members are in some way similar to each other [35].

Received: 9 March 2016 Accepted: 9 August 2016 Published online: 20 August 2016

## References

1. Zachary WW (1977) An information flow model for conflict and fission in small groups. *J Anthropol Res* 33(4):452-473
2. Newman ME (2012) Communities, modules and large-scale structure in networks. *Nat Phys* 8(1):25-31
3. Fortunato S (2010) Community detection in graphs. *Phys Rep* 486(3):75-174
4. Lancichinetti A, Fortunato S, Radicchi F (2008) Benchmark graphs for testing community detection algorithms. *Phys Rev E* 78(4):046110

5. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E (2008) Fast unfolding of communities in large networks. *J Stat Mech Theory Exp* 2008(10):P10008
6. Bullmore E, Sporns O (2009) Complex brain networks: graph theoretical analysis of structural and functional systems. *Nat Rev Neurosci* 2009(10):186-198
7. Shibata N, Kajikawa Y, Takeda Y, Matsushima K (2008) Detecting emerging research fronts based on topological measures in citation networks of scientific publications. *Technovation* 28(11):758-775
8. Herrera M, Roberts DC, Gulbahce N (2010) Mapping the evolution of scientific fields. *PLoS ONE* 5(5):e10355
9. Rosvall M, Bergstrom CT (2010) Mapping change in large networks. *PLoS ONE* 5(1):e8694
10. Chen P, Redner S (2010) Community structure of the physical review citation network. *J Informetr* 4(3):278-290
11. Hric D, Darst RK, Fortunato S (2014) Community detection in networks: structural communities versus ground truth. *Phys Rev E* 90(6):062805
12. Leskovec J, Adamic LA, Huberman BA (2007) The dynamics of viral marketing. *ACM Trans Web* 1(1):5
13. Backstrom L, Huttenlocher D, Kleinberg J, Lan X (2006) Group formation in large social networks: membership, growth, and evolution. In: *Proceedings of the 12th ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, New York, pp 44-54
14. Mislove A, Marcon M, Gummadi KP, Druschel P, Bhattacharjee B (2007) Measurement and analysis of online social networks. In: *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*. ACM, New York, pp 29-42
15. Yang J, Leskovec J (2015) Defining and evaluating network communities based on ground-truth. *Knowl Inf Syst* 42(1):181-213
16. Palchykov V, Kaski K, Kertész J, Barabási A-L, Dunbar RI (2012) Sex differences in intimate relationships. *Sci Rep* 2:370
17. Kovanen L, Kaski K, Kertész J, Saramäki J (2013) Temporal motifs reveal homophily, gender-specific patterns, and group talk in call sequences. *Proc Natl Acad Sci USA* 110(45):18070-18075
18. Expert P, Evans TS, Blondel VD, Lambiotte R (2011) Uncovering space-independent communities in spatial networks. *Proc Natl Acad Sci USA* 108(19):7663-7668
19. Bothorel C, Cruz JD, Magnani M, Micenkova B (2015) Clustering attributed graphs: models, measures and methods. *Netw Sci* 3(3):408-444
20. Newman MEJ, Clauset A (2016) Structure and inference in annotated networks. *Nat Commun* 7:11863
21. Girvan M, Newman ME (2002) Community structure in social and biological networks. *Proc Natl Acad Sci USA* 99(12):7821-7826
22. Waltman L, Eck NJ (2012) A new methodology for constructing a publication-level classification system of science. *J Am Soc Inf Sci Technol* 63(12):2378-2392
23. Boyack KW, Klavans R (2010) Co-citation analysis, bibliographic coupling, and direct citation: which citation approach represents the research front most accurately? *J Am Soc Inf Sci Technol* 61(12):2389-2404
24. Boyack KW, Newman D, Duhon RJ, Klavans R, Patek M, Biberstine JR, Schijvenaars B, Skupin A, Ma N, Börner K (2011) Clustering more than two million biomedical publications: comparing the accuracies of nine text-based similarity approaches. *PLoS ONE* 6(3):e18029
25. Glenisson P, Glänzel W, Janssens F, De Moor B (2005) Combining full text and bibliometric information in mapping scientific disciplines. *Inf Process Manag* 41(6):1548-1572
26. An electronic archive and distribution server for research articles. <http://arxiv.org>
27. Prokofyev R, Demartini G, Boyarsky A, Ruchayskiy O, Cudré-Mauroux P (2013) Ontology-based word sense disambiguation for scientific literature. In: *European conference on information retrieval*. Springer, Berlin, pp 594-605.
28. Jones KS (1973) Index term weighting. *Inf Storage Retr* 9(11):619-633. doi:10.1016/0020-0271(73)90043-0
29. Newman ME, Girvan M (2004) Finding and evaluating community structure in networks. *Phys Rev E* 69(2):026113
30. Blei DM, Ng AY, Jordan MI (2003) Latent Dirichlet allocation. *J Mach Learn Res* 3:993-1022
31. Barber MJ (2007) Modularity and community detection in bipartite networks. *Phys Rev E* 76(6):066102
32. Guimerà R, Sales-Pardo M, Amaral LAN (2007) Module identification in bipartite and directed networks. *Phys Rev E* 76(3):036102
33. Larremore DB, Clauset A, Jacobs AZ (2014) Efficiently inferring community structure in bipartite networks. *Phys Rev E* 90(1):012805
34. Meilă M (2007) Comparing clusterings - an information based distance. *J Multivar Anal* 98(5):873-895
35. Jacob EK (2004) Classification and categorization: a difference that makes a difference

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](http://springeropen.com)

---