# Identifying Hidden Communities of Interest with Topic-based Networks: A Case Study of the Community of Philosophers of Science (1930-2017)

Christophe Malaterre
Département de philosophie & Centre interuniversitaire de recherche sur la science et la technologie, Université du Québec à Montréal (UQAM), 455 Boulevard René-Lévesque Est, Montréal (QC) H3C 3P8, Canada
malaterre.christophe@uqam.ca

Francis Lareau
Département d'informatique, Université du Québec à Montréal (UQAM), 201 Avenue Président-Kennedy, Montréal (QC) H2X 3Y7, Canada
lareau.francis@courrier.uqam.ca

## Abstract

*Scientific networks are often investigated by means of citation analyses. Yet, interpretation of such networks in terms of semantic (and often disciplinary) content heavily depends on supplementary knowledge, notably about author research specialties. Similar situations arise more generally in many types of social networks whose semantic interpretation relies on supplementary information. Here, author community networks are inferred from a topic model which provides direct insights into the semantic specificities of the identified "hidden communities of interest" (HCoI). Using a philosophy of science corpus of full-text articles (N=16,917), we investigate its underlying communities by measuring topic profile correlations between authors. A diachronic perspective is built by modeling the research networks over different time-periods and mapping genealogical relationships between communities. The results show a marked increase in philosophy of science communities over time and trace the progressive appearance of the specialization areas that structure the field today.*

**Keywords:** hidden communities of interest, hidden colleges, social networks, text-mining, topic-modeling, philosophy of science

## 1. Introduction

Analysis of the network of relationships between different actors is well-known to provide a wealth of insights, be it about the identification of key members and relationships, or the structure of the network itself and its evolution over time. This is what social network analysis (SNA) is all about, having developed into a research field of its own (Molontay & Nagy, 2019). SNA has notably been applied to social networks formed by scientists, leading to research on latent social structures or "hidden colleges" (Crane, 1969) and to a broader understanding of the dynamics of science, including the role of networks of scientists and institutions over such issues as problem selection, discovery, collaboration or even career dynamics (Fortunato et al., 2018). Yet, before being submitted to analysis, social networks first need to be built, which presupposes access to data about actors and their relationships. One frequently used source of such data comes from social media. In science, academic social media have been mined to identify researchers profiles, collaborations and trajectories (Tang et al., 2008; Kong et al., 2019). Another major source is citation data: citation analyses, and more generally bibliometric approaches, have been shown to provide reliable insights into the structure of scientific social networks (Small, 1999; Boyack et al., 2005). This has led to numerous studies applied to different domains of science, including even the field of SNA itself (Molontay & Nagy, 2019).

However, interpretation of such social networks usually requires supplementary information from which to infer the meaning of relationships between actors or the specificity of communities. For instance, in the case of scientific networks, this can be information about author research specialties which is often obtained by examining key publications or metadata (e.g. keywords) from other sources (Raimbault et al., 2016). With such approaches, networks are built from relational data (links between nodes) and the meaning of the resulting communities (set of nodes which may share similar content) is inferred with the help of supplementary data. To address this issue more systematically, some have proposed to develop specific topic models that could incorporate author-related data (Steyvers et al., 2004), notably in the case of social media and directional networks (McCallum et al., 2007; Pathak et al., 2008), or in the case of co-authorship data (Zhou et al., 2006; Zhang et al., 2007). Others have proposed to further develop community detection algorithms so as to include not just topological information but also prior constraining data on nodes, as in semi-supervised community detection algorithms and graph neural network approaches (Yang et al., 2014; Ye et al., 2018).

HICSS

The present work tackles the question of social networks in an extreme context where semantic information about actors is abundant but relational data are scarce or even inexistant. Such situations could be seen as even precluding the very notion of social network. Yet, underlying communities of actors can still be identified based on their shared semantic content. Such communities are so to speak "hidden communities of interest" (HCoI), that is to say groups of actors sharing similar semantic contents but whose social relationships with one another may be unknown. HCoI's reflect the existence of underlying latent social networks whose study can nevertheless be pursued to gain insights, for instance, into their structure and evolution (even if usual social network metrics such as centrality should receive a different interpretation).

Here, we investigate the concrete structure of HCoI's in a specific scientific context by applying a combination of topic-modeling and community detection approaches. We also show how diachronic analyses can be performed to reveal the temporal evolution of these HCoI's.

We do so on an academic corpus of research articles in the philosophy of science. Text-mining approaches have already been applied to broadly map the discipline of philosophy (Buckner et al., 2011), and so have SNA and citation analysis approaches (Noichl, 2021). SNA has also been applied to the more specialized field of the philosophy of science with a view to identifying its key journals (Wray, 2010), studying its relationship with the domain of the history of science (Weingart, 2015), and assessing its impact onto scientific publications (Khelfaoui et al., 2021). In parallel, topic-based approaches have been used to mine the key journals of the field and provide a more detailed view of the research themes of the field and their evolution in the 20[th] century (Malaterre & Lareau, 2022b). Yet to date, no investigation of communities of interest based on shared research themes found in publications has ever been conducted. This is what the present work is about.

More specifically, the proposed approach starts by fitting an LDA topic-model to this philosophy of science corpus, thereby resulting in topic probability distributions for all full-text articles. Having split the corpus into four broad time-periods, these topic probability distributions were averaged out per author and per time-period, depending on authors contribution to each article. This resulted in author topic profiles for each time-period. Correlation analyses between author topic profiles then led to the construction of author correlation networks for each time-period, which were submitted to Louvain community detection. In turn, the topic profile of each community was quantified by averaging out their author topic profiles. These community topic profiles provide immediate insights into the semantic specificities of each community. Furthermore, measuring pairwise distances between community topic profiles across time-periods provides a means to understand the diachronic evolution of communities and their genealogies. In what follows, we first describe in more details the data and methods (Section 2). We then present the results, notably the networks of communities that were detected and their temporal evolution (Section 3). These findings are then discussed (Section 4).

## 2. Data and methods

For this study, we used a corpus of full-text articles from eight major philosophy of science journals that had been assembled in (Malaterre & Lareau, 2022a). The corpus spans from 1930 (the first issue of the earliest published journal) to 2017 and includes 16,917 research articles from the *British Journal for the Philosophy of Science*, the *European Journal for Philosophy of Science*, *Erkenntnis*, *International Studies in the Philosophy of Science*, the *Journal for General Philosophy of Science*, *Philosophy of Science*, *Studies in History and Philosophy of Science Part A* and *Synthese*. The corpus was cleaned and preprocessed in a standard way. Only nouns, verbs, adverbs and adjectives were kept following part-of-speech (POS) tagging and lemmatization (TreeTagger package (Schmid, 1994) with Penn Tree-Bank tag sets (Marcus et al., 1993)) and words occurring in fewer than 50 sentences in the corpus were removed.

Following (Malaterre & Lareau, 2022a), topic-modeling was carried out with the well-known Latent Dirichlet Allocation (LDA) algorithm, following (Blei et al., 2003) and (Griffiths & Steyvers, 2004), with a number of topics $K$ set to 25 (this number was chosen as a compromise between an optimal coherence measure following (Röder et al., 2015) and upon manual inspection of models with lower and higher $K$ values). This therefore resulted in 25 probability distributions over the corpus terms (each probability distribution considered to represent a topic), and the probability distributions of these topics in each one of the 16,917 articles. Inspection of the most probable terms within each topic and of selected text excerpts made it possible to interpret and label all topics. For ease of handling, topics were also grouped into categories based on their correlation within corpus documents and Louvain community detection performed on the graph of topic correlations in Gephi (Bastian et al., 2009). These categories were interpreted based on expert knowledge of the field.

In parallel, all author metadata were manually checked and disambiguated, ensuring similar spellings were used throughout the corpus (especially in the early decades). All authors ($N$=8009) were assigned publication weights based on their respective number of articles (coauthored articles were evenly split). Four main time-

periods of 21 years each were then defined (1930-1951, 1952-1973, 1974-1995, 1996-2017) and article topic distributions were averaged out per author for each one of these periods. This step resulted in topic profiles for each author based on their publications during any given time-period (in other words: for each time-period, author-specific probability distributions over the 25 topics).

For each time-period, Pearson correlations among these author topic profiles were calculated. Correlation networks were built in Gephi (Bastian et al., 2009), using Louvain community detection (with default parameters). To reduce noise, only authors with weighted publication above 2 were retained (thereby filtering out "transient authors", see section 3.1), and a correlation threshold was set to 0.6 (this resulted in keeping all significant author communities connected to the network main component across all 4 time-periods while removing clutter). Topic profiles (i.e., topic probability distributions) were then calculated for each community by averaging out their author topic profiles. To get further insights into the genealogy of communities over time, we calculated the Hellinger distances between community topic profiles across time-periods with the Gensim implementation (Rehurek & Sojka, 2010) and focused on the closest pairings, thus generating a diachronic picture of the evolution of philosopher communities and their main research themes.

## 3. Results

### 3.1. Authors and articles in the corpus

The corpus covers the content of eight of the most significant philosophy of science journals in English language. Of course, philosophy of science is published in numerous other venues, including general philosophy journals (e.g., *Mind*), disciplinary focused journals (e.g., *Biology and Philosophy*), science journals or books. It is also published in many non-English languages. Nevertheless, the corpus we have assembled includes the most authoritative journals of the field. It also includes the journals that started the field in the 1930s and are still flagship journals today. These are good reasons for accepting the corpus as offering a representative perspective of the discipline. Out of the 16,917 articles written by 8,009 authors, about 70% come from the three earliest published journals: *Erkenntnis, Philosophy of Science* and *Synthese* (Table 1). The other five journals are later comers to the field, the most recent one being the *European Journal for the Philosophy of Science* founded in 2011.

Over the past eight decades, the volume of articles has significantly increased, from 1,575 for the 1930-1951 period to 8,300 for the 1996-2017 period, which is

a 5.3-fold increase (Fig. 1). Meanwhile the number of authors has incurred an 8-fold increase. Knowing that the number of articles per author has roughly remained constant throughout all four periods at about 2, the increase in authors denotes an increase in co-authorship. Indeed, the number of multi-authored articles has increased 4-fold, from 4% in the first period to 16% in the last. Though this share is significantly lower than in the sciences where single-authored articles are now virtually non-existent (e.g. in ecology, Barlow et al., 2018), or even in some areas of the humanities (e.g. in economics, Kuld & O'Hagan, 2018), multiple-authorship has been steadily rising in the philosophy of science.

**Table 1. Corpus journals**

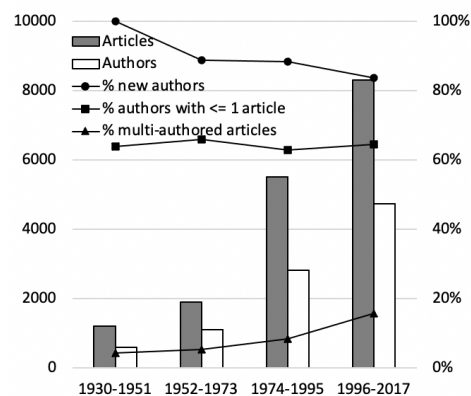| Journals (alphabetic order) | Publication periods | Total |
|---|---|---|
| *British Journal for the Philosophy of Science* | 1950-present | 1,862 |
| *Erkenntnis* | 1930-1940; 1975-present | 2,127 |
| *European Journal of Philosophy of Science* | 2011-present | 156 |
| *International Studies in the Philosophy of Science* | 1986-present | 560 |
| *Journal for General Philosophy of Science* | 1970-present | 929 |
| *Philosophy of Science* | 1934-present | 4,605 |
| *Studies in History and Philosophy of Science A* | 1970-present | 1,421 |
| *Synthese* | 1936-1939, 1946-1949, 1955-present | 5,257 |
| Total | 1930-2017 | 16,917 |



**Figure 1. Authors and articles**

Note that the proportion of authors who only publish once (or "transients", see (Crane, 1969)) is relatively stable at about 65%. This not significantly different from what is observed elsewhere, a partial explanation being the share of doctoral students and post-doctoral researchers (e.g. in synthetic biology, Raimbault et

al., 2016). Moreover, while the number of new authors (from one period to another) is above 80%, this proportion has been decreasing over time. This means that, although many authors in any given period were not present in the previous period, new authors now tend to represent a smaller share of authors compared to what it used to be.

## 3.2. Topics

The results of the topic-modeling provide a high-level perspective on the main research themes of the philosophy of science (Malaterre & Lareau, 2022a). These topics can be analyzed by examining their top-words (as in Table 2) as well as sample texts. Group A of topics denote research questions that are characteristic of the philosophy of language and logic. Group B includes topics in epistemology and theory of knowledge (including questions about realism), while group C relates more specifically to induction, confirmation, and the use of probabilities. Group D is about rational decision and game theory. Topics in the philosophy of biology and the neurosciences are found in group E. Group F includes a set of traditional topics which concern the process of scientific explanation, the nature of causation and the status of natural kinds. Topics in the philosophy of physics are in group G, with thermodynamics, electromagnetism, chemistry in one topic, and relativity and quantum theory in the other. Finally, group H gathers topics which are characterized by a more historical or social discourse. These include research themes in the history of science and in the history of philosophy, but also investigations on the social dimensions of science.

## 3.3. Author communities through time

Over the past eight decades, the philosophy of science has significantly grown in terms both of research domains and authors. As can be seen in Fig. 2A, the field comprised a handful of communities in the 1930s-40s whose topic profiles are depicted in Fig. 2B. A clearly identifiable cluster (1a) consists in the community of the logical positivists and members of the Vienna circle (e.g., Neurath, Reichenbach, Carnap, Hempel), distinctively focused on philosophy of language and logic (Fig. 2A; note that author name size and node size are proportional to weighted number of publications). As is well-known, the subsequent development of the philosophy of science owes much to these authors. The core of the network consists of two closely interconnected communities, with, on the one hand (1c), a group of researchers somehow at the border between philosophy and other humanities (e.g., history, anthropology, economics, psychology), and on the other (1d) authors engaging in more traditional metaphysics or ontology (e.g., realism, subjectivity etc.). Though engaging with science, these two groups remained much anchored to a classical philosophical discourse. A distinct and much smaller community (1b), somehow at the fringe of the core of the network, consists in philosophers focusing more on physics, and discussing issues related to matter, energy, or physical theories (e.g., electromagnetism or quantum mechanics; note the presence of Malisoff, founder of *Philosophy of Science*).

### Table 2. Topics and keywords

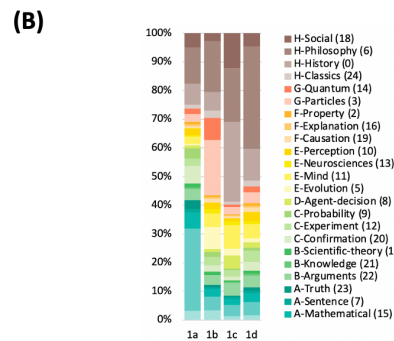| Topic name | Top-10 words |
|---|---|
| A-Formal | set; function; relation; define; definition; structure; order; model; theory; class |
| A-Language | language; sentence; term; meaning; concept; use; statement; logical; mean; word |
| A-Mathematical | mathematical; mathematics; number; proof; axiom; geometry; theory; object; point; line |
| A-Sentence | sentence; context; use; say; reference; content; name; true; semantic; speaker |
| A-Truth | logic; truth; true; proposition; sentence; logical; formula; follow; rule; world |
| B-Arguments | argument; claim; say; question; make; view; reason; fact; case; point |
| B-Knowledge | belief; knowledge; epistemic; believe; know; case; evidence; reason; justification; true |
| B-Scientific-theory | theory; scientific; theoretical; empirical; realism; realist; truth; science; true; claim |
| C-Confirmation | law; hypothesis; statement; evidence; theory; condition; inductive; problem; confirmation; fact |
| C-Experiment | datum; experiment; value; use; test; result; experimental; model; hypothesis; method |
| C-Probability | probability; measure; value; give; chance; case; function; distribution; degree; frequency |
| D-Agent-decision | agent; action; decision; game; choice; act; utility; strategy; moral; preference |
| E-Evolution | selection; population; organism; evolutionary; gene; biological; individual; group; evolution; specie |
| E-Mind | behavior; state; mental; action; psychological; human; function; psychology; person; child |
| E-Neurosciences | system; information; process; cognitive; level; mechanism; state; representation; structure; function |
| E-Perception | object; experience; perception; see; color; perceptual; visual; content; red; image |
| F-Causation | causal; cause; event; effect; causation; condition; case; variable; time; occur |
| F-Explanation | model; explanation; explain; account; explanatory; phenomenon; use; case; system; provide |
| F-Property | property; world; object; physical; relation; kind; entity; part; identity; exist |
| G-Particles | theory; energy; law; particle; electron; atom; physical; physic; chemical; system |
| G-Quantum | time; state; space; quantum; system; theory; particle; physical; field; point |
| H-Classics | motion; body; force; newton; law; galileo; earth; move; light; time |
| H-History | work; time; man; history; new; year; make; life; century; write |
| H-Philosophy | world; nature; knowledge; concept; experience; kant; sense; thing; idea; reality |
| H-Social | science; scientific; social; research; scientist; philosophy; knowledge; problem; history; practice |



Figure 2. (A) Author communities and (B) their topic profiles 1930-1951

The discipline developed substantially in the 1950s throughout the early 1970s, with an increasing number of interconnected communities. As can be seen on Fig. 3A, community (2a) includes logicians and philosophers of mathematics with a distinctively formal vocabulary. Philosophers of language constitute a separate community (2b). Occupying a central position in the network (2c), a community targeting specific issues related to confirmation and the status of scientific theories (e.g., induction, verifiability, corroboration, or refutation; note the presence of Popper). A small and peripheric group of researchers (2d) consists of the nascent community of philosophers of biology, with a notable focus on evolutionary theory. On the opposite, philosophers of physics constitute a larger community (2e), addressing a diversity of epistemic issues related for instance to relativity theory or quantum mechanics. In continuity with the previous period, a distinctive community is constituted by authors at the border with traditional philosophy (2f), while a nearby community appears to address more sociological aspects of science (2g).
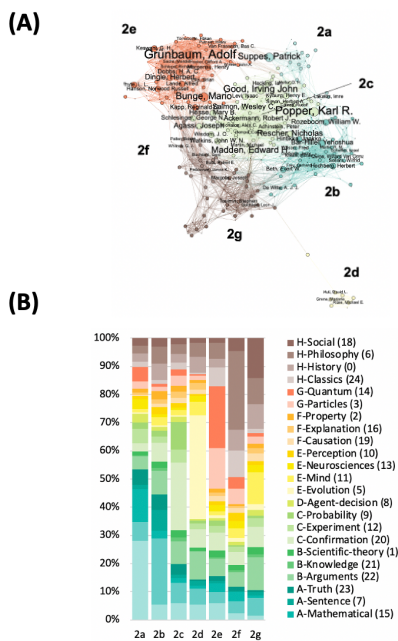
centered on semantics and the philosophy of language (3b; note the presence of Hintikka known for his work on formal epistemic logic and of game semantics for logic). A specific community consists of authors addressing epistemology and theory of knowledge questions (3c). Questions about confirmation and the status of scientific theories characterize a community somehow at the center of the network (3d). Note the appearance of a specific community focused on probability theory and its relevance for science and knowledge (3e). Another new community consists of authors interested in decision and game theories, and their applications in science (3f).



**Figure 3. (A) Author communities and (B) their topic profiles 1952-1973**

In the 1970s throughout the 1990s, the philosophy of science continued to grow in terms of authors but also in terms of topic communities (Fig. 4). Community (3a) consists of logicians and philosophers of mathematics, somehow in continuity with a second community more



**Figure 4. (A) Author communities and (B) their topic profiles 1974-1995**

The community of philosophers of biology remains at the margin of the rest of philosophy of science but has significantly grown in size (3g; note the presence of Sober). A novel community has appeared around researchers more specifically targeting the philosophy of mind and the neurosciences (3h). Yet another community of philosophers distinctively focuses on causation (3i). Two communities are characterized by topics related to the philosophy of physics: a first one with a characteristic focus on quantum mechanics and relativity (3k), and

a second smaller one more oriented towards thermodynamics, chemistry, or electromagnetism (3j). A relatively diffuse community gathers philosophers who tend to have a more traditional philosophical standpoint (3l). Finally, a large community consists in a diverse set of authors who tend to target some social dimensions in science (3m).

The trend towards an increase in terms of number of authors and a specialization of discursive topics continued in the 1990s throughout the 2010s (Fig. 5). A community of philosophers of language (4a) can be seen quite tightly connected to a second community of philosophers of logic (including modal and intuitionistic logic) notably interested in notions of truth (4b). A nearby community consists of epistemologists, philosophers specializing in theory of knowledge (4c).
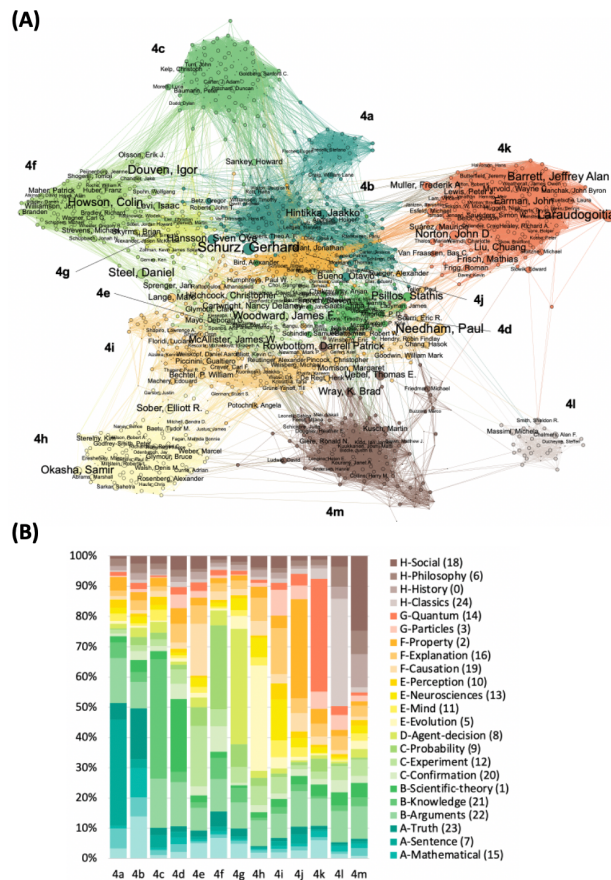


**Figure 5. (A) Author communities and (B) their topic profiles 1996-2017**

Towards the center of the network, a community focuses on the status of scientific theories notably with respect to realist and anti-realist stances (4d). A nearby

and more diffuse community gathers authors interested in topics that relate to data, experiments, and modeling, but also somehow to causation (4e). The community of philosophers of probability, which had appeared in the previous period has grown in size and individuated itself (4f). A community of researchers somehow bridging philosophers of probability and of logic consists of authors focusing on game theory and various aspects of rational choice theory (4g). The community of philosophers of biology (4h) has significantly grown and is somehow more integrated with the rest of the network, notably with the community of philosophers of the neurosciences and others interested in scientific explanation (4i). At the center of the network lies a community generally interested in ontology (4j, addressing issues about properties or kinds among others). A large group of philosophers of physics constitutes a relatively well distinct community that tends to focus on relativity and quantum mechanics, with related issues such as the structure of space-time (4k). A noticeably distinct category appears to mobilize classical philosophical works in their discussion of science (4l). Finally, a community of authors focuses on the social dimensions of science and various aspects of the practice of science (4m).

### 3.4. Retracing community genealogies

Measuring the distances between the topic profiles of any two communities from two different periods provides insights on the transformation of HCoI's into one another through time: the shorter the distances, the closer the communities in terms of their topic interests (Fig. 6).

Transitioning from the first period (1930-1951) to the second (1952-1973) (Fig. 6A), one sees a reasonably good filiation between communities focused on philosophy of language and logic (1a to 2b). However, the other three communities tend to consolidate into one, and possibly a second (1b-d to 2f-g): the early philosophy of physics (1b) tends to bifurcate into a community still centered on similar physics-related topics (2e) and another community closer to traditional philosophy (2f), the latter being in the continuity of authors engaging in more metaphysics or ontology (1d). Note how the communities 2f and 2g appear to be relatively close to all the communities of the previous period, indicating a reconfiguration of authors and their topics of interest. Given the increase in the number of communities, this also denotes a form of marginalization of what once constituted the core of the philosophy of science. Note how the philosophy of biology community of the second period (2d) shows little continuity with previous communities, indicating the emergence of a novel HCoI.

The transition from the second period (1952-1973) to the third (1974-1995) also shows an increase in the

number of communities, yet filiations tend to be stronger, indicating a form of stabilization of research communities with novel themes still emerging (Fig 6B).
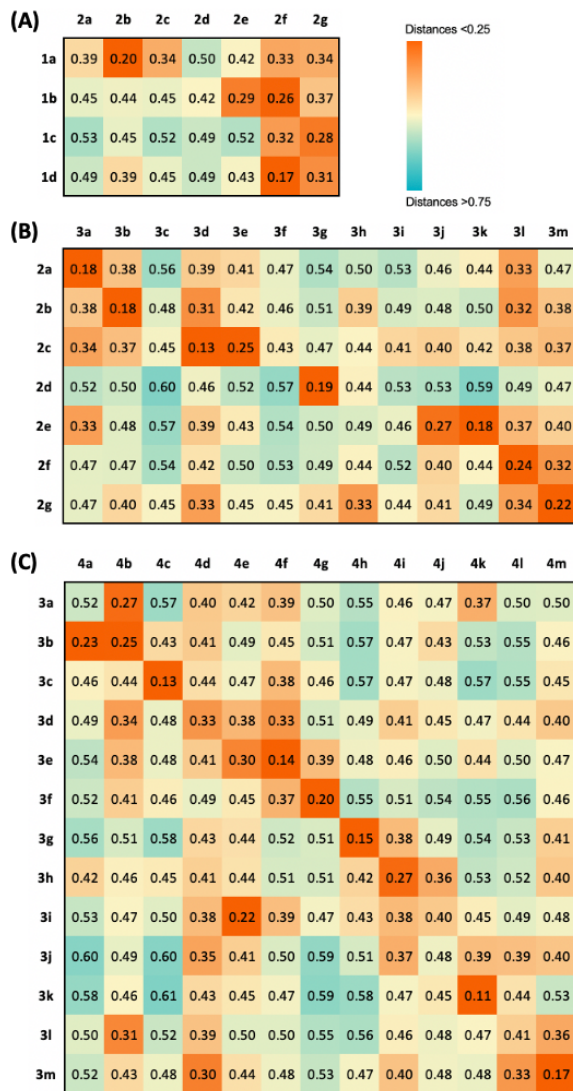
**(A)**

| | 2a | 2b | 2c | 2d | 2e | 2f | 2g |
|---|---|---|---|---|---|---|---|
| 1a | 0.39 | 0.20 | 0.34 | 0.50 | 0.42 | 0.33 | 0.34 |
| 1b | 0.45 | 0.44 | 0.45 | 0.42 | 0.29 | 0.26 | 0.37 |
| 1c | 0.53 | 0.45 | 0.52 | 0.49 | 0.52 | 0.32 | 0.28 |
| 1d | 0.49 | 0.39 | 0.45 | 0.49 | 0.43 | 0.17 | 0.31 |

Distances <0.25

Distances >0.75

**(B)**

| | 3a | 3b | 3c | 3d | 3e | 3f | 3g | 3h | 3i | 3j | 3k | 3l | 3m |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2a | 0.18 | 0.38 | 0.56 | 0.39 | 0.41 | 0.47 | 0.54 | 0.50 | 0.53 | 0.46 | 0.44 | 0.33 | 0.47 |
| 2b | 0.38 | 0.18 | 0.48 | 0.31 | 0.42 | 0.46 | 0.51 | 0.39 | 0.49 | 0.48 | 0.50 | 0.32 | 0.38 |
| 2c | 0.34 | 0.37 | 0.45 | 0.13 | 0.25 | 0.43 | 0.47 | 0.44 | 0.41 | 0.40 | 0.42 | 0.38 | 0.37 |
| 2d | 0.52 | 0.50 | 0.60 | 0.46 | 0.52 | 0.57 | 0.19 | 0.44 | 0.53 | 0.53 | 0.59 | 0.49 | 0.47 |
| 2e | 0.33 | 0.48 | 0.57 | 0.39 | 0.43 | 0.54 | 0.50 | 0.49 | 0.46 | 0.27 | 0.18 | 0.37 | 0.40 |
| 2f | 0.47 | 0.47 | 0.54 | 0.42 | 0.50 | 0.53 | 0.49 | 0.44 | 0.52 | 0.40 | 0.44 | 0.24 | 0.32 |
| 2g | 0.47 | 0.40 | 0.45 | 0.33 | 0.45 | 0.45 | 0.41 | 0.33 | 0.44 | 0.41 | 0.49 | 0.34 | 0.22 |

**(C)**

| | 4a | 4b | 4c | 4d | 4e | 4f | 4g | 4h | 4i | 4j | 4k | 4l | 4m |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3a | 0.52 | 0.27 | 0.57 | 0.40 | 0.42 | 0.39 | 0.50 | 0.55 | 0.46 | 0.47 | 0.37 | 0.50 | 0.50 |
| 3b | 0.23 | 0.25 | 0.43 | 0.41 | 0.49 | 0.45 | 0.51 | 0.57 | 0.47 | 0.43 | 0.53 | 0.55 | 0.46 |
| 3c | 0.46 | 0.44 | 0.13 | 0.44 | 0.47 | 0.38 | 0.46 | 0.57 | 0.47 | 0.48 | 0.57 | 0.55 | 0.45 |
| 3d | 0.49 | 0.34 | 0.48 | 0.33 | 0.38 | 0.33 | 0.51 | 0.49 | 0.41 | 0.45 | 0.47 | 0.44 | 0.40 |
| 3e | 0.54 | 0.38 | 0.48 | 0.41 | 0.30 | 0.14 | 0.39 | 0.48 | 0.46 | 0.50 | 0.44 | 0.50 | 0.47 |
| 3f | 0.52 | 0.41 | 0.46 | 0.49 | 0.45 | 0.37 | 0.20 | 0.55 | 0.51 | 0.54 | 0.55 | 0.56 | 0.46 |
| 3g | 0.56 | 0.51 | 0.58 | 0.43 | 0.44 | 0.52 | 0.51 | 0.15 | 0.38 | 0.49 | 0.54 | 0.53 | 0.41 |
| 3h | 0.42 | 0.46 | 0.45 | 0.41 | 0.44 | 0.51 | 0.51 | 0.42 | 0.27 | 0.36 | 0.53 | 0.52 | 0.40 |
| 3i | 0.53 | 0.47 | 0.50 | 0.38 | 0.22 | 0.39 | 0.47 | 0.43 | 0.38 | 0.40 | 0.45 | 0.49 | 0.48 |
| 3j | 0.60 | 0.49 | 0.60 | 0.35 | 0.41 | 0.50 | 0.59 | 0.51 | 0.37 | 0.48 | 0.39 | 0.39 | 0.40 |
| 3k | 0.58 | 0.46 | 0.61 | 0.43 | 0.45 | 0.47 | 0.59 | 0.58 | 0.47 | 0.45 | 0.11 | 0.44 | 0.53 |
| 3l | 0.50 | 0.31 | 0.52 | 0.39 | 0.50 | 0.50 | 0.55 | 0.56 | 0.46 | 0.48 | 0.47 | 0.41 | 0.36 |
| 3m | 0.52 | 0.43 | 0.48 | 0.30 | 0.44 | 0.48 | 0.53 | 0.47 | 0.40 | 0.48 | 0.48 | 0.33 | 0.17 |

**Figure 6. Community distances (Hellinger distances between community topic probability distributions)**

Philosophy of language and logic communities map well onto one another (2a-b to 3a-b). The community about confirmation and scientific theories (2c) persisted into 3d, while giving rise to a distinct community focused on probability theory and its relevance for knowledge (3e). The philosophy of biology community also persisted as a well identified set of authors and topics (2d to 3g). The community of philosophers of physics (2e) appears to have grown and split into one community more focused on relativity and quantum theory (3k) and another on the rest of physics (3j). The two socio-historico-philosophical communities (2f-g) somehow persisted (3l-m), though one notes a relative proximity of the latter communities to many of the communities of the previous period, indicating multiple reconfigurations. Four novel communities appeared in the 1970s-80s without any clear filiation from communities of the previous period: a community focusing on knowledge theory (3c), another exploring game theory and rational choice (3f), yet another on philosophy of mind and the neurosciences (3h), and finally a community distinctively focusing on causation (3i).

The number of communities stabilized during the last decades of the 20th century. The transition from the third (1974-1995) to the fourth period (1996-2017) shows a relatively good continuity (Fig. 6C). Communities of philosophers of language and logic slightly reorganized themselves depending on topic alignments but remained stable as a group (3a-b to 4a-b). Epistemologists persisted as a specific community, while gaining in momentum and autonomy (3c to 4c). The community focusing on probability theory and knowledge also persisted (3e to 4f), as well the communities on game theory (3f to 4g), on the philosophy of biology (3g to 4h), on the philosophy of mind and the neurosciences (3h to 4i), on the philosophy of relativity and quantum theory (3k to 4k), and on the social dimensions of knowledge (3m to 4m). On the other hand, some communities tend to have somehow dissolved into several. This is notably the case for the community on confirmation and scientific theories (3d) denoting a detachment from these topics in the 1990s-2000s. Similarly, the community focusing on chemistry, electromagnetism, or thermodynamics (3j) has somehow dissolved in subsequent decades, as well as the one which was centered on more traditional philosophical issues (3l). Finally, philosophers focusing on causation (3i) appear to have joined a broader community also interested in data, experiments, and modeling (4e).

## 4. Discussion

The identification of what we called "hidden communities of interest" (HCoI) in the philosophy of science and their mapping through time reveals the semantic proximity of certain authors with one another and how these latent intellectual groups distinctively evolved. The results highlight semantic reconfigurations in the field, with communities dissolving into others (e.g., the community of confirmation and scientific theories of the 1970s-1990s), while novel and well-delineated communities appeared (e.g., the communities of philosophers of biology in the early 1970s, of epistemologists and of philosophers interested in game theory in the 1980s). Overall, the evolution of author communities shows a phase of growth and diversification as the

number of authors (and interests) increased, followed by a later phase of stabilization characterized by a form of intellectual entrenchment of larger and usually well delineated communities. These results concur with known episodes of the field, for instance the role of logical positivism in the constitution of the philosophy of science in the early 20th century (Giere & Richardson, 1996) or the emergence of a philosophy of biology in the 1970s as can be reconstituted by examining dedicated anthologies (Sober, 2006; Rosenberg & Arp, 2009), while providing a richer mapping of the diversity of author groups, their relative proximity and their temporal evolution. The results also provide a quantitative basis for what is often considered implicit knowledge of the field. In sum, they offer a high-level account of the recent history of the field and its actors – a history which remains to be written in all its details. The findings certainly also raise a number of questions, among which the evolution of the relatively fuzzy and malleable communities characterized by a form of socio-historico-philosophical discourse, or the reasons accounting for the relative isolation of specific communities such as epistemologists, or the still the factors explaining changes in the degree of specificity of still other communities such as philosophers interested in causation or in confirmation. As such, the community networks provide exploratory tools with which to formulate further questions about the structure of the field, its history and the role and place of specific authors.

As we have seen, the findings result from a combination of topic-modeling and community detection approaches. The main objective of these approaches is to identify hidden communities of interest, that is to say groups of actors sharing similar semantic contents but whose social relationships with one another may be unknown. The methods thereby make it possible to identify communities of actors on the basis of their semantic content in the absence of known social connections. They also make it possible to assess the relative topic proximity of these communities. The resulting networks differ from social networks as usually construed. Indeed, social networks depict actual relationships between agents. In science, these relationships often result from bibliometric assessments, for instance measures of citations, co-citations or co-authorships (Boyack et al., 2005). Here, the author networks are based on the similarity of their topic profiles (averaged from their respective publications). Consequently, it is a priori possible that two authors from the same community may not appear nearby one another in a bibliometric-based social network (for instance if they do not cite the same literature though discussing the same topics, or if they do not cite each other). And conversely, it is a priori possible that two closely related authors in a bibliometric network may happen to be in different topic communities

(for instance if they do not work on the same research questions though citing the same references). In practice, we doubt such cases would be frequent, though this is something that would need to be specifically measured in further work (this could be done by combining both types of approaches, for instance on a full-text corpus with lists of references). One caveat of the approach is that typical analyses of social networks (e.g., centrality) do not have the same meaning in topic-based communities: here, a most central author simply is one whose topic profile is the most similar to those of others; yet such an author does not need to be the most central in terms of citations. Similarly, to make significant authors stand out, especially compared to transients, we opted for network representations in which node size was proportional to the number of publications in the period. Yet publication volume is definitely not always indicative of citation impact (though conversely, it could lead us to reconsider whose works to read and cite). Again, combining both bibliometric and topic-based approaches should make it possible to alleviate such concerns. On the other hand, one major advantage of the HCoI approach is to provide a topic chart for each community. Whereas bibliometric approaches must rely on supplementary investigations about author profiles in order to make sense of the observed networks, the topic-based network approach makes it possible to immediately understand the specific identity of each community in terms of discursive topics.

Analyzing the corpus in specific time-periods resulted in a sequence of several topic-based author networks, making it possible to assess the temporal dynamics at play, notably increases in size or in number of communities. Measuring the pairwise distances between community topic profiles provides a mapping of communities from one period to the next. This additional layer of information gives insights on community filiation and their overall genealogy. In the case of the present corpus, this approach made it possible to understand the appearance (and sometimes disappearance) of specific communities that constitute the discipline of the philosophy of science. It should also be possible to sort communities into different types depending on their diachronic fate. Several cases could be distinguished, notably: one-to-one relationships (cases of clear filiation), one-to-several (bifurcation), one-to-none (dissolution), several-to-several (re-configuration), several-to-one (consolidation), none-to-one (emergence). Further work could look at automatically classifying communities along such lines.

## 5. Conclusion

Combining topic-modeling and community detection methods makes it possible to uncover hidden

communities of interest (HCoI) and map their proximity in terms of semantic content both synchronically (through correlation networks) and diachronically (over time-periods). Using a working corpus of 16,917 full-text research articles written by 8,009 philosophers of science from the 1930s up to the 2010s, this approach revealed how these authors constituted usually well delineated HCoI's characterized by specific topic profiles. The results notably show how the discipline of philosophy of science grew and diversified in terms of research themes and communities over the past eight decades. The approach makes it possible to gain insights into author-based communities, notably their semantic content in the form of directly interpretable topic profiles, but also their relative proximity and temporal evolution. When data about actual social interactions are not available but textual data are, mapping such HCoI networks can still provide relevant insights about groups of social actors sharing similar interests. In cases where both textual and social data are available, HCoIs analyses could also provide a complementary content-based perspective compared to usual social network analyses.

# 6. References

Barlow, J., Stephens, P. A., Bode, M., Cadotte, M. W., Lucas, K., Newton, E., Nuñez, M. A., & Pettorelli, N. (2018). On the extinction of the single-authored paper: The causes and consequences of increasingly collaborative applied ecological research. *Journal of Applied Ecology*, *55*(1), 1–4. https://doi.org/10.1111/1365-2664.13040

Bastian, M., Heymann, S., & Jacomy, M. (2009). Gephi: An Open Source Software for Exploring and Manipulating Networks. *International AAAI Conference on Weblogs and Social Media*.

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, *3*(Jan), 993–1022.

Boyack, K. W., Klavans, R., & Börner, K. (2005). Mapping the backbone of science. *Scientometrics*, *64*(3), 351–374. https://doi.org/10.1007/s11192-005-0255-6

Buckner, C., Niepert, M., & Allen, C. (2011). From encyclopedia to ontology: Toward dynamic representation of the discipline of philosophy. *Synthese*, *182*(2,), 205–233.

Crane, D. (1969). Social Structure in a Group of Scientists: A Test of the "Invisible College" Hypothesis. *American Sociological Review*, *34*(3), 335. https://doi.org/10.2307/2092499

Fortunato, S., Bergstrom, C. T., Börner, K., Evans, J. A., Helbing, D., Milojević, S., Petersen, A. M., Radicchi, F., Sinatra, R., Uzzi, B., Vespignani, A., Waltman, L., Wang, D., & Barabási, A.-L. (2018). Science of science. *Science*, *359*(6379), eaao0185. https://doi.org/10.1126/science.aao0185

Giere, R. N., & Richardson, A. W. (Eds.). (1996). *Origins of logical empiricism*. University of Minnesota Press.

Griffiths, T. L., & Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National Academy of Sciences*, *101*(suppl 1), 5228–5235. https://doi.org/10.1073/pnas.0307752101

Khelfaoui, M., Gingras, Y., Lemoine, M., & Pradeu, T. (2021). The visibility of philosophy of science in the sciences, 1980–2018. *Synthese*. https://doi.org/10.1007/s11229-021-03067-x

Kong, X., Shi, Y., Yu, S., Liu, J., & Xia, F. (2019). Academic social networks: Modeling, analysis, mining and applications. *Journal of Network and Computer Applications*, *132*, 86–103. https://doi.org/10.1016/j.jnca.2019.01.029

Kuld, L., & O'Hagan, J. (2018). Rise of multi-authored papers in economics: Demise of the 'lone star' and why? *Scientometrics*, *114*(3), 1207–1225. https://doi.org/10.1007/s11192-017-2588-3

Malaterre, C., & Lareau, F. (2022a). On the use of machine translation and topic-modeling to analyze non-parallel multilingual corpora: A case study in the history of philosophy of science. *Zenodo.Org*. https://doi.org/10.5281/zenodo.6484582

Malaterre, C., & Lareau, F. (2022b). The early days of contemporary philosophy of science: Novel insights from machine translation and topic-modeling of non-parallel multilingual corpora. *Synthese*, *200*(3), 242. https://doi.org/10.1007/s11229-022-03722-x

Marcus, M. P., Marcinkiewicz, M. A., & Santorini, B. (1993). Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, *19*(2), 313–330. https://doi.org/10.21236/ADA273556

McCallum, A., Wang, X., & Corrada-Emmanuel, A. (2007). Topic and Role Discovery in Social Networks with Experiments on Enron and Academic Email. *Journal of Artificial Intelligence Research*, *30*, 249–272. https://doi.org/10.1613/jair.2229

Molontay, R., & Nagy, M. (2019). Two Decades of Network Science as seen through the co-authorship network of network scientists. *Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, 578–583. https://doi.org/10.1145/3341161.3343685

Noichl, M. (2021). Modeling the structure of recent philosophy. *Synthese*, *198*, 5089–5100. https://doi.org/10.1007/s11229-019-02390-8

Pathak, N., DeLong, C., & Banerjee, A. (2008). *Social Topic Models for Community Extraction*. 10.

Raimbault, B., Cointet, J.-P., & Joly, P.-B. (2016). Mapping the Emergence of Synthetic Biology. *PLOS ONE*, *11*(9), e0161522. https://doi.org/10.1371/journal.pone.0161522

Rehurek, R., & Sojka, P. (2010). Software Framework for Topic Modelling with Large Corpora. *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, 45–50.

Röder, M., Both, A., & Hinneburg, A. (2015). Exploring the Space of Topic Coherence Measures. *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining - WSDM '15*, 399–408. https://doi.org/10.1145/2684822.2685324

Rosenberg, A., & Arp, R. (Eds.). (2009). *Philosophy of*

*Biology: An Anthology.* Wiley-Blackwell.

Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. *Proceedings of International Conference on New Methods in Language Processing*, 44–49.

Small, H. (1999). Visualizing science by citation mapping. *Journal of the American Society for Information Science*, *50*(9), 799–813. https://doi.org/10.1002/(SICI)1097-4571(1999)50:9<799::AID-ASI9>3.0.CO;2-G

Sober, E. (Ed.). (2006). *Conceptual issues in evolutionary biology* (3rd ed). MIT Press.

Steyvers, M., Smyth, P., Rosen-Zvi, M., & Griffiths, T. (2004). Probabilistic author-topic models for information discovery. *Proceedings of the 2004 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining  - KDD '04*, 306. https://doi.org/10.1145/1014052.1014087

Tang, J., Zhang, J., Yao, L., Li, J., Zhang, L., & Su, Z. (2008). ArnetMiner: Extraction and mining of academic social networks. *Proceeding of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD 08*, 990. https://doi.org/10.1145/1401890.1402008

Weingart, S. B. (2015). Finding the History and Philosophy of Science. *Erkenntnis*, *80*(1), 201–213. https://doi.org/10.1007/s10670-014-9621-1

Wray, K. B. (2010). Philosophy of Science: What are the Key Journals in the Field? *Erkenntnis*, *72*(3), 423–430. https://doi.org/10.1007/s10670-010-9214-6

Yang, L., Cao, X., Jin, D., Wang, X., & Meng, D. (2014). A unified semi-supervised community detection framework using latent space graph regularization. *IEEE Transactions on Cybernetics*, *45*(11), 2585–2598.

Ye, F., Chen, C., & Zheng, Z. (2018). Deep autoencoder-like nonnegative matrix factorization for community detection. *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 1393–1402.

Zhang, H., Qiu, B., Giles, C. L., Foley, H. C., & Yen, J. (2007). An LDA-based community structure discovery approach for large-scale social networks. *2007 IEEE Intelligence and Security Informatics*, 200–207.

Zhou, D., Ji, X., Zha, H., & Giles, C. L. (2006). Topic evolution and social interactions: How authors effect research. *Proceedings of the 15th ACM International Conference on Information and Knowledge Management*, 248–257.