# The evolution of scientific literature as metastable knowledge states

**Sai Dileep Koneru**[1], **David Rench McCauley**[2], **Michael C. Smith**[2], **David Guarrera**[2], **Jenn Robinson**[2], **and Sarah Rajtmajer**[1,*]

[1]The Pennsylvania State University, University Park, PA, USA
[2]Ernst & Young, McLean, VA, USA
[*]smr48@psu.edu

## ABSTRACT

The problem of identifying common concepts in the sciences and deciding when new ideas have emerged is an open one. Metascience researchers have sought to formalize principles underlying stages in the life-cycle of scientific research, determine how knowledge is transferred between scientists and stakeholders, and understand how new ideas are generated and take hold. Here, we model the state of scientific knowledge immediately preceding new directions of research as a metastable state and the creation of new concepts as combinatorial innovation. We find that, through the combined use of natural language clustering and citation graph analysis, we can predict the evolution of ideas over time and thus connect a single scientific article to past and future concepts in a way that goes beyond traditional citation and reference connections.

## Introduction

Early work in metascience can be traced back at least half a century,[1] although it has been only in the last decade or so that a robust literature has been seeded exploring co-authorship networks, citation networks, topical networks and similar static and one-dimensional representations of complex interactions amongst researchers and their work. Much of this has been powered by the increased availability of digital data on scientific processes, improvements in information retrieval, network science, machine learning, and computational power, allowing researchers to derive meaningful insights. A substantial subset of this literature has focused on quantifying and predicting success in publishing – how we should measure success, who will have it, and what factors contribute to having it. Seminal work has focused on modeling citation patterns for papers[2] and researchers[3], with more recent work setting out to explain hot streaks in researchers' career trajectories[4], unique patterns of productivity and collaboration amongst the scientific elite[5], and even the role of luck in driving scientific success[6,7]. We are also seeing the emergence of metascience as a social movement[8], catalyzed by the last decade's reproducibility crisis[9], aiming to describe and evaluate science at a macro scale in order to diagnose biases in research practice[10,11], highlight flaws in publication processes[12], understand how researchers select new work to pursue[13,14], identify opportunities for increased efficiency (e.g., automated hypothesis generation[15]), and forecast emergence of research topics[16,17].

Prior work on the evolution of research can be broadly viewed in three categories based on method: network-based, language-based, and hybrid methods using both networks and language. Language-based methods commonly use topic models such as Latent Dirichlet Allocation (LDA) and predict changes in topics[18,19]. Other language-based approaches include tracking usage of keywords[20], analyzing linguistic context[16], and modeling topics sequentially[17]. Studies using network-based methods usually use citation networks and community detection algorithms such as topological clustering methods[21] or clique percolation methods[22] to identify emergence of new fields, while other network approaches include usage of temporal[23], multiplex[24] networks, projections of citation networks such as co-authorship[25,26]. Hybrid usage of both language- and network-based methods to predict the evolution of scientific fields includes keyword-generated networks used to predict changes in topics[27] or approaches that mostly rely on network analysis, applying linguistic techniques such as LDA for explanatory labels only[28]. Still others[29] have used LDA and co-occurrence networks of topics to study changes in knowledge-based systems. However, to the best of our knowledge, these hybrid methods do not incorporate state-of-the-art language embeddings, nor do they incorporate insights from both the language models and the citation network. It is our hypothesis that the only way to truly capture the evolution of ideas and knowledge in the literature is through the integration of network and linguistic techniques.
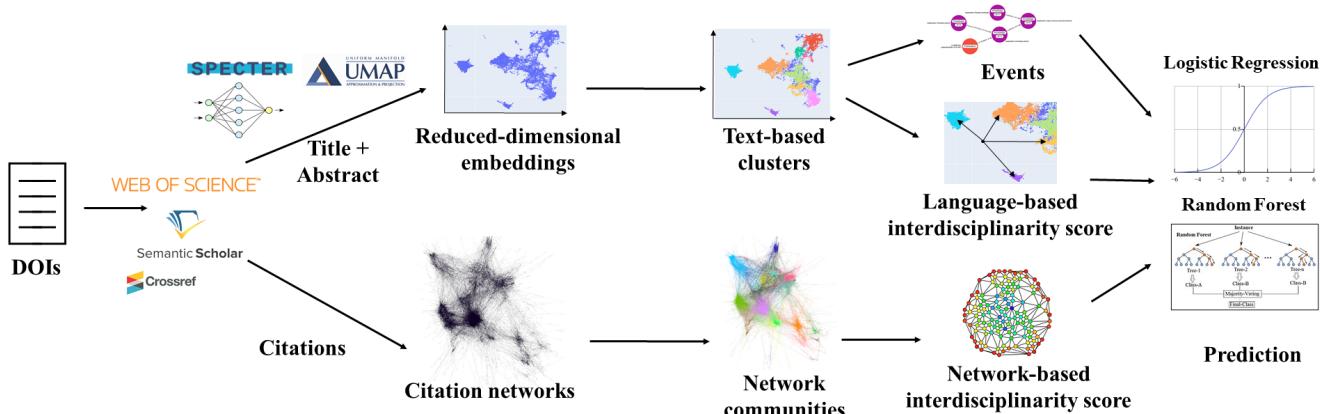
A premise of the work discussed in this paper is that neither citation networks alone (or derivatives thereof) nor purely language-driven models of the scientific corpus can explain the evolution of fields and the emergence of new ideas. We show that these two frameworks capture overlapping but distinct and complementary aspects of dynamics in scientific research. We use pre-trained neural network models[30] to generate vectorized representations of the literature while separately leveraging citation network measures (e.g., betweenness centrality), combining these two inputs to build predictive models of topical

evolution. The intuition behind the mechanisms explored herein is that scientific disciplines can be described at a high level by aggregation of related ideas. When a discipline is beginning to show signs of fracture or change via the emergence or synthesis of new ideas, we model this moment borrowing from physics the concept of *metastability*: a state easily perturbed into a new state. We suggest that measures of interdisciplinarity may be indicators of this transition and thus useful predictors of change in the scientific ecosystem.

Recent efforts have elevated the role of interdisciplinarity in scientific practice[31–34]. Prior work has shown interdisciplinarity to have an effect on innovation and research impact[11,35]. Calls for collaboration across disciplines are prominent throughout research institutions and funding agencies[1] but some have argued that the promises of interdisciplinarity are overstated and misplaced[36]. The bibliometric community has offered a data-driven framing for interdisciplinary studies, e.g., defining interdisciplinarity as a process of integrating different bodies of knowledge[37,38].

The definition of interdisciplinarity varies broadly in the literature, with different definitions capturing different aspects of this concept[39], and can be broadly classified into two groups: subject-based and network-based definitions[39]. Subject-based metrics rely on multi-classification systems to calculate interdisciplinarity, leaning on pre-defined subject categories, e.g., from the Web of Science (WoS)[40]. These approaches generally are imposed at the journal level, focusing on the distribution of subject categories, e.g., percentage of references cited by publications in journals outside a journal of interest's category[41,42]. In some cases, metrics are borrowed from other fields, such as the Gini index from economics and Shannon entropy from information theory, to quantify diversity[43]); these are also based on pre-made categories. Alternatively, network-based interdisciplinarity metrics are often assessed based on the location of a publication in a citation network[44], with centrality measures frequently being the focus. For example, betweenness centrality, which is independent of third-party categorization, was one of the first metrics used in this way[44,45] and has likewise been used to predict future network trends[46,47].

To study knowledge evolution in the scientific literature, we: (1) develop methods that utilize transformers-based language models and unsupervised clustering to track the evolution of ideas over time; (2) quantify interdisciplinarity using complementary text- and citation-based metrics; and (3) explore the utility of metastability, measured through interdisciplinarity, as a predictor of scientific evolutionary events (Fig. 1).



**Figure 1. Data analysis workflows.** (Top) **Text-based analysis.** Title and abstract are concatenated and input to a language embedding model, then dimensionally reduced and fed into a clustering algorithm; clusters of embedded papers are then used for event modeling and interdisciplinarity scoring. (Bottom) **Citation-based analysis.** Citation information is used to create undirected citation graphs; the Louvain algorithm is used to identify network communities and betweenness centrality is used for interdisciplinarity scoring. Interdisciplinary metrics are jointly used to predict disciplinary evolution.

# Dataset

Our dataset contains detailed records of 19,177 scientific papers published in the years 2011 through 2018, with 2300 to 2500 papers for each year, representing a substantial stratified random sample of papers published in 62 prominent journals from the following disciplines as strata: Criminology; Economics and Finance; Education; Health; Management; Marketing and Organizational Behavior; Political Science; Psychology; Public Administration; and Sociology.[2] Metadata for these were collected using the Web of Science as a primary source. Digital Object Identifiers (DOIs) were used to merge WoS records with Semantic Scholar (S2) records[49,50] for completeness of metadata coverage and author name disambiguation. When DOIs

---

[1]See, e.g., the U.S. National Science Foundation's Growing Convergence Research program: https://www.nsf.gov/od/oia/growing-convergence-research/index.jsp.

[2]The dataset was collected in conjunction with DARPA's SCORE program. For a complete listing of journals see[48].

were not available from WoS, we used Crossref[51] to fill in missing DOIs for more complete record linking between the Web of Science and Semantic Scholar. For citation network analyses, we also included all papers referenced by these papers. Our complete dataset includes records of 839,096 papers and about 1.45 million citations.

## Methods

We use parallel workflows to model dynamics in bibliometric data – one based on text and one based on citation networks (Fig. 1). For each we derive a measure of interdisciplinarity useful for prediction of knowledge evolution and we describe our explanatory and predictive experiments to evaluate our measures.

### SPECTER-based topic modeling

We use language-embedding-based topic modeling to identify topics within our corpus for a given year. To do so, we extract embeddings for each publication in our dataset using the concatenated title and abstract as an input to SPECTER (Scientific Paper Embeddings using Citation-informed TransformERs)[30], a model for generating document-level embeddings of scientific documents via pre-training on scientific papers and their citation graphs.[3] SPECTER embeddings have been shown to outperform competitive baselines on benchmark document-level tasks such as citation prediction, document classification and recommendation[30].

To identify disciplines and subdisciplines, we use an unsupervised, non-parametric, hierarchical clustering algorithm, Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN)[53]. Specifically, we soft-cluster SPECTER embeddings to reflect that papers may belong to multiple (sub)disciplines with different probabilities. As the performance of HDBSCAN generally reduces as the dimensionality of input data increases, we use UMAP[54] to reduce the dimensionality of SPECTER embeddings prior to clustering with HDBSCAN. We use multi-objective Bayesian hyperparameter tuning[55] for the UMAP-HDBSCAN pipeline to balance five evaluative criteria related to balancing inter- vs. intra-cluster density, number of clusters, and persistence of clusters over multiple runs of the algorithm. The successfully clustered papers are considered "strong members" of that cluster.

We refer to papers that cannot be confidently assigned by the clustering algorithm as "weak members". We assign each weak member to the cluster with which it has the highest semantic similarity. Downstream analyses are reported with and without inclusion of weak members. We consider this distinction because we suggest that weak members represent research which is significantly different (and potential truly innovative) relative to existing disciplines, and as such can help explain shifts in the trajectories of fields.[4]

For each cluster, we generate representative keyphrases using a procedure similar to the KeyBERT library[56], with modifications (e.g., more performant aggregation of embeddings from large numbers of documents belonging to the same cluster). Deriving keyphrases provides explanatory power for clusters and adds more nuanced understanding of the clusters than other commonly used approaches to grouping knowledge products, e.g., WoS categories. Clusters identified in our dataset for the year 2011 and their corresponding keyphrases are shown in Figure 2. Our approach identifies a total of 371 clusters over the dataset, i.e., years 2011 through 2018.
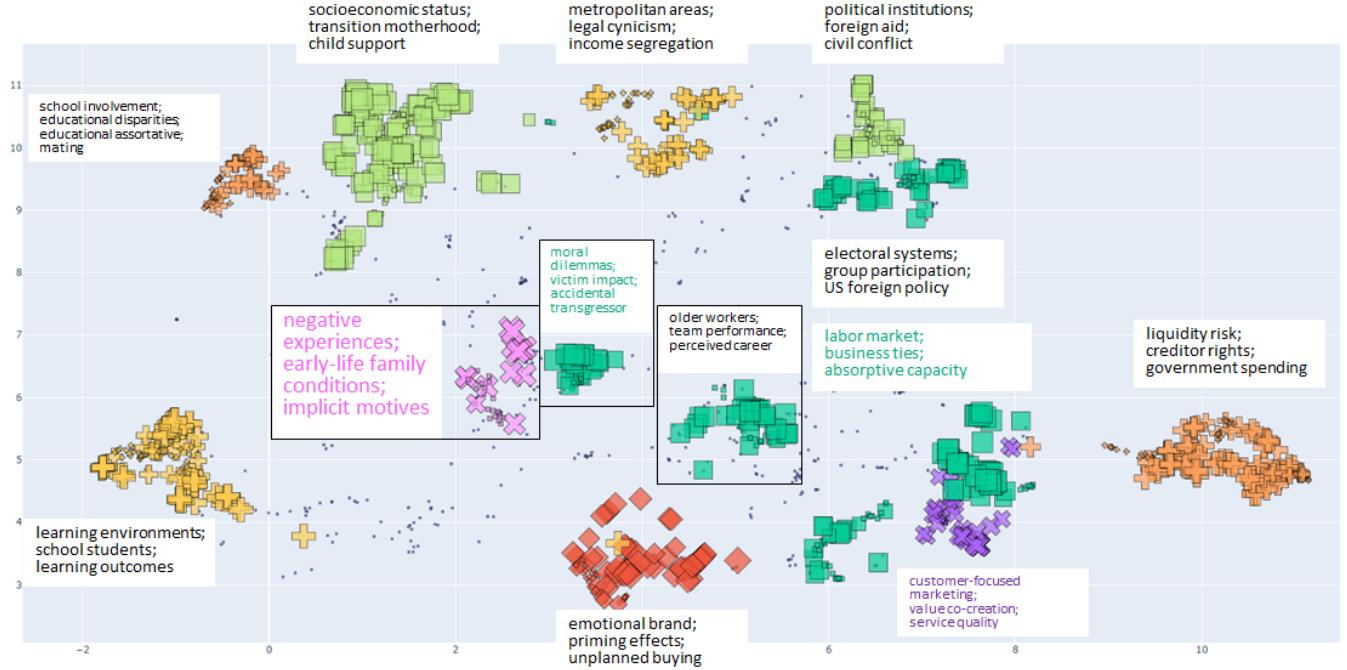
### Citation graphs and communities

Per common practice, our citation-based analysis considers the citation network wherein nodes in the graph represent papers in our dataset and undirected edges represent citation relationships. We detect communities in this network using the Louvain community detection algorithm[57]. Commonly-used, it maximizes modularity of the network, namely the expected value of inter- vs intra-community edges[58]. Specifically, for a given time window/year of interest $t$ we consider the subgraph $G(t)$ containing only papers published in year $t$ and earlier, as well as their references. This approach allows us to make predictions for past papers without fear that future papers citing them will cause information leakage into the dataset (e.g. a model trying to predict the evolution of an idea tied to a paper from 2017 should not have access to information about papers from 2018 citing it during model training). An example of the community structure discovered via the Louvain method is shown in Figure 3.

### Quantifying interdisciplinarity

**Language-based interdisciplinarity:** Our text-based interdisciplinarity (ID) metric scores each publication based on its soft clustering membership probabilities (i.e. the probability of a publication belonging to each possible cluster identified by standard or "hard" clustering), considering only strong member publications. It does so by assuming that one representation of

---

[3]Specifically, we use the `huggingface` implementation[52] of the pre-trained SPECTER model.

[4]HDBSCAN refers to these non-confident assignments as noise; however, we expect these not to be noise in the traditional sense (e.g., an outlier or data worthy of discarding as it provides no analytical value) but instead to potentially add value as extremely novel research.

**Figure 2.** UMAP projection of publications in the year 2011 colored by HDBSCAN-generated cluster labels with corresponding cluster-level keyphrases. Each cluster plotted here contains at least 2.5% of total papers for the year and the size of each point is proportional to that publication's language-based interdisciplinarity score. Small blue points represent weak members. Note that most clusters shown are well-separated and not homogeneous in shape, suggesting that UMAP is doing a good job of dimensionally reducing the feature space in such a way that it is reasonably straightforward to partition and that a variable-density-based clustering algorithm, such as HDBSCAN, is well-suited to identifying clusters in such a dataset.

interdisciplinarity is the diversity of language pulled from different fields. This metric is calculated using Equation 1 which considers the spread in its cluster assignment probabilities. Formally:

$$\text{ID}_{text} = \frac{N}{N-1} \left(1 - P_{wm} - \max\left(P_{cluster}\right)\right)\left(1 - \sigma_p\right) \tag{1}$$

where $N$ is the total number of clusters in the dataset, $P_{cluster}$ is the probability of the paper belonging to a cluster, $P_{wm}$ is the probability of the paper being a weak member of any cluster, and $\sigma_p$ the standard deviation of $P_{cluster}$ over all clusters. This formulation is more intuitive when extreme cases are considered. For example, consider a corpus with 9 clusters for the year of interest. Consider a paper that sits very clearly within a single well-defined scientific discipline, i.e., $max(P_{cluster}) = 1$ for a single cluster (consequently, $P_{wm} = 0$). The interdisciplinarity score for that paper would be $ID_{text} = 0.0$. Alternatively, imagine a paper with membership probabilities that are equivalent for all clusters, with the same probability that it may be a weak member, i.e., $P_{wm} = P_{cluster,i} = 0.1$ for $N = 9$. This would result in $ID_{text} = 0.9$, reflecting that the paper belongs to a wide array of disciplines/clusters equally, but also there is some chance that it may be a weak member – which can also be interpreted as a global uncertainty in the membership probabilities – thus keeping it from achieving a score of 1.0.
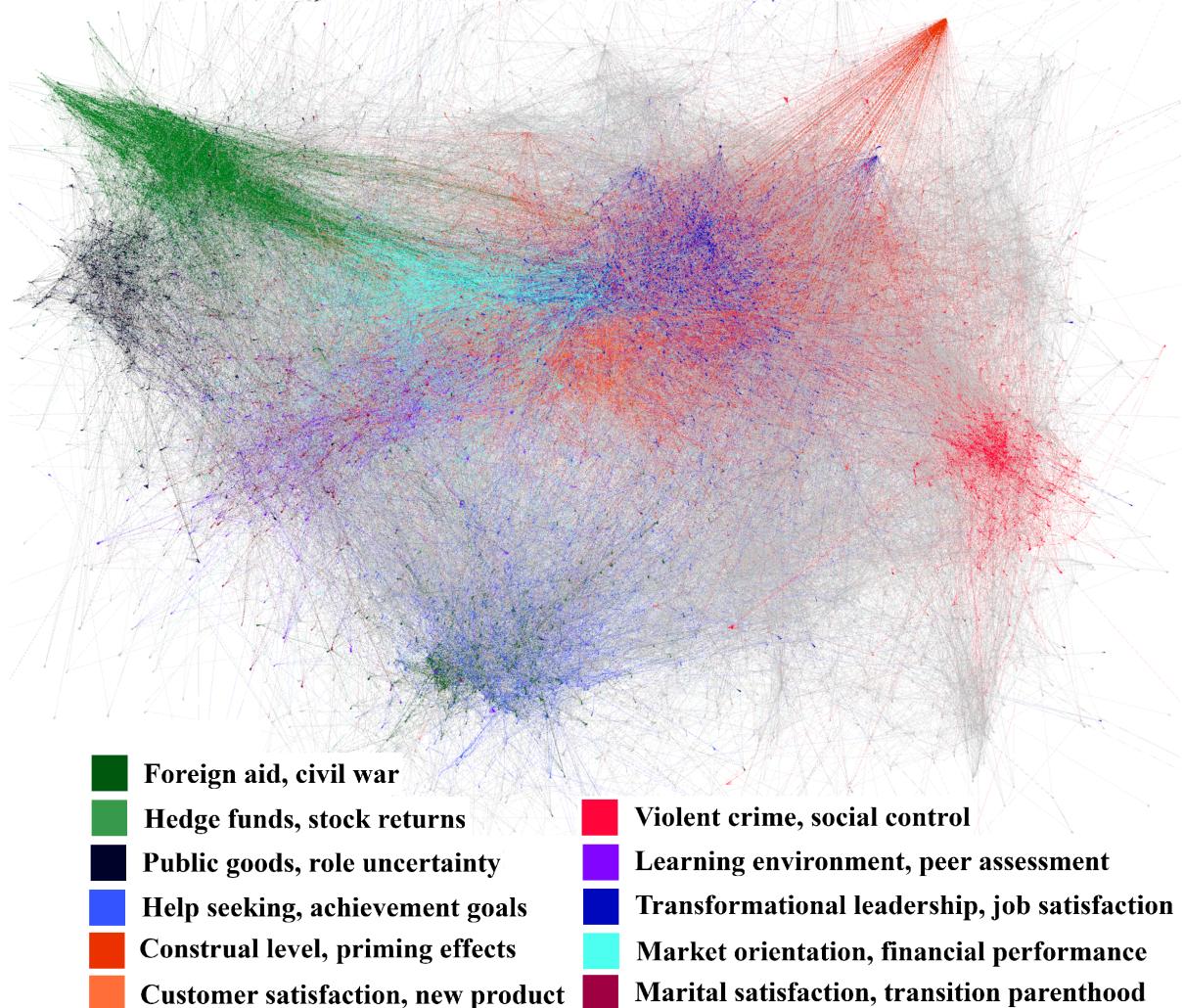
**Citation-based interdisciplinarity:** We use betweenness centrality for each publication in the network as an interdisciplinarity metric, with higher centrality generally indicating higher interdisciplinarity, as has been done in previous literature[59]. As we do for community detection, we use time-windowed subgraphs for centrality measurement. Betweenness centrality is lightly modified for use as an ID metric, normalized on a [0,1] scale. For paper $i$ in publication year $t$:

$$\text{ID}_{network} = centrality_{t,i}/max(\{centrality_t\}) \tag{2}$$

where $\{centrality_t\}$ is the set of all centrality values for papers published in calendar year $t$.

## Text-based dynamic event modeling

We identify and track critical knowledge evolution events borrowing from the literature tracking communities in dynamic social networks[60]. Specifically, representative embeddings for each cluster are calculated using the element-wise mean of embeddings of the papers in each cluster, and clusters are compared across consecutive years by calculating the pairwise cosine similarity of

**Foreign aid, civil war**

**Hedge funds, stock returns**

**Violent crime, social control**

**Public goods, role uncertainty**

**Learning environment, peer assessment**

**Help seeking, achievement goals**

**Transformational leadership, job satisfaction**

**Construal level, priming effects**

**Market orientation, financial performance**

**Customer satisfaction, new product**

**Marital satisfaction, transition parenthood**

**Figure 3.** An exemplary snapshot of the dense network and communities found by the Louvain community detection algorithm for the year 2011. Communities comprising less than 2.5% of total papers for the year are colored grey. Note that a clear community structure can be observed for this graph-only approach much like it was for the language-only clustering presented earlier.

the embeddings of each $[C_t, C_{t+1}]$ pair of clusters in years $t$ and $t+1$[60]. We then link a cluster with its best-matching cluster(s) in the consecutive time step if the cosine similarity is above 0.95. We employ the following taxonomy[60]:

- A *birth* event is identified at time $t$ when a cluster at time $t$ has no matching cluster(s) at time $t-1$.

- A *death* event is identified at time $t$ when a cluster at time $t$ has no matching cluster(s) at time $t+1$.

- Multiple clusters have *merged* at time $t$ when one cluster at time $t$ matches to two or more clusters at time $t-1$.

- Multiple clusters have *split* at time $t$ when one or more clusters at time $t$ match to a single cluster at time $t-1$.

- A *continuation* event is observed when one cluster at time $t$ is matched to exactly one cluster at time $t+1$.

We group these events into two types for subsequent analyses: (1) *dynamic* – split or merge and (2) *stable* – continuity or death.[5] Figure 4 gives a notional example of merge and continuation events. We note that events may occur in combination; e.g., a cluster may split into two, and those two clusters may simultaneously merge with two other clusters.

### Event-tracking and prediction

We hypothesize that interdisciplinarity scores and cluster size are indicators of metastability and therefore can be used to predict cluster evolution, i.e., *dynamic* vs. *stable* events, as an endogenous and target variable. In particular, for each language cluster $C_t$ at time $t$, we use as exogenous model inputs: cluster-wise mean language-based interdisciplinarity score (which does include weak member papers); mean citation-based interdisciplinarity score for weak and strong members, treated as separate features in order to discern if there is any difference in predictive power considering weak members; and number of weak and strong member papers in the cluster.

To choose the most powerful features and test their predictive power (and thus value for further analyses), we use multinomial logistic regression and a Random Forest classifier with a binary target $\vec{y}$ representing if a dynamic event type (split or merge) is observed for a cluster at time $t+1$ as shown in Equation 3.

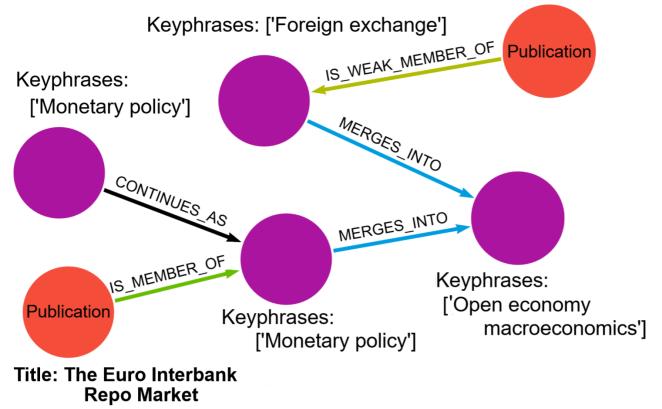$$\vec{y} = \begin{bmatrix} split/merge \\ continuation/death \end{bmatrix} \quad (3)$$

We use the entire dataset with multinomial logistic regression for explanatory power. For the random forest, we use cluster events in the period 2011-17 for training and the year of 2018 for testing, resulting in roughly an 86%/14% train/test split by cluster count with 275 events for training (split/merge: 136; continuation/death: 139) and 43 testing events (split/merge: 21, continuation/death: 22). Using the above input features and event types in year $t+1$, we fit a random forest with 100 trees using the default hyperparameter values from the `scikit-learn` python library[61].



**Figure 4.** Notional continuation and merge events showing weak (significantly different from existing clusters) and strong members (high confidence in its membership) of each cluster.

## Results

In the following, we first show that language and network frameworks capture different information by comparing the overlap between clusters identified using text and citation-based communities. We then further investigate the nature of the information provided by both frameworks by discussing how these representations, when considered together, not only serve to predict the evolution of disciplines and sub-fields but are equally important when doing so.

### Comparing clusters and communities suggests valuable incomplete overlap

Figure 3 gives a snapshot of network communities in 2011; comparison with Figure 2 illustrates differences in grouping across the two approaches. In general, the Louvain algorithm detects communities in the citation network at a finer resolution than our text-based clustering. For reference, Figure 5 shows the number of clusters and communities in our dataset, in addition to a measure of overlap between the two that we describe below. The number of network communities generally decreases over time, reflecting a more integrated citation graph emerging amongst the papers in our sample.

---

[5]Not only does treating splits and merges as a single class emerge from our metastability mental model but, given that they often co-occur, this treatment creates non-overlapping classes. We disregard birth events at present since a birth event has no preceding data from which to build a model and is unrelated to the concept of combinatorial innovation being described by metastability.

As both our language- and citation-based frameworks are unsupervised, to compare them we need to identify clusters with one another across frameworks. For this, we measure pairwise Jaccard similarity between clusters and communities, effectively looking at the fraction of shared publications between every language cluster and every network community relative to their total number of member papers. If the similarity between a cluster and a community is above 0.1 then we consider them similar. This threshold-based method (and the 0.1 threshold specifically) has been used in the literature for tracking clusters and communities over time[60,62] and performs well across a variety of synthetic graphs. Going back to Figure 5, the inset shows the percentage of language clusters with similar (Jaccard similarity > 0.1) network communities. It can be seen that while there is an overlap between the communities and clusters, the overlap is not complete, which suggests that each approach adds unique insight.



**Figure 5.** Plot with number of clusters/communities identified by text-based (brown) and networks-based (blue) frameworks with inset plot showing percentage of language clusters associated with at least one network-derived community. Note that overlap values are consistently below 100% but well above 0%, suggesting unique and complementary insights added by each.

**Illustrating knowledge evolution events**

To illustrate the types of knowledge events we identify and track in this work, let us consider an example from our dataset. Figure 6 shows the evolution of a full chain of language cluster evolutionary events over the period 2011 through 2018. Every cluster in this chain has 'Business and Finance' and 'Economics' as the most common WoS categories among member papers. In contrast, the keyphrases generated via our language clustering approach (shown in the Figure) reflect a much greater resolution, including phrases like "income hedging" and "intangible capital". This chain starts with a 2011 cluster that appears related to the (then recent) U.S. housing market crisis and Great Recession. There is a strong focus on work discussing corporate governance and government spending. This focus on organizational-level finance and economics mostly continues through 2017, with only a few deviations that are more focused on overall market trends. This is epitomized by the representative paper for one of the 2016 clusters, focused on European banking. Then something happens in 2018: topics appear to shift substantially from organizational/macroeconomic concepts to research focused on individual-level spending, finance, and decision-making, as can be seen both from the keyphrases representing those linguistic clusters, as well as from the representative 2018 paper focused on accounting for consumer behaviors in investing. It is interesting to note that this timing corresponds with Richard Thaler's 2017 Nobel Prize in Economics, awarded for contributions to behavioural economics.
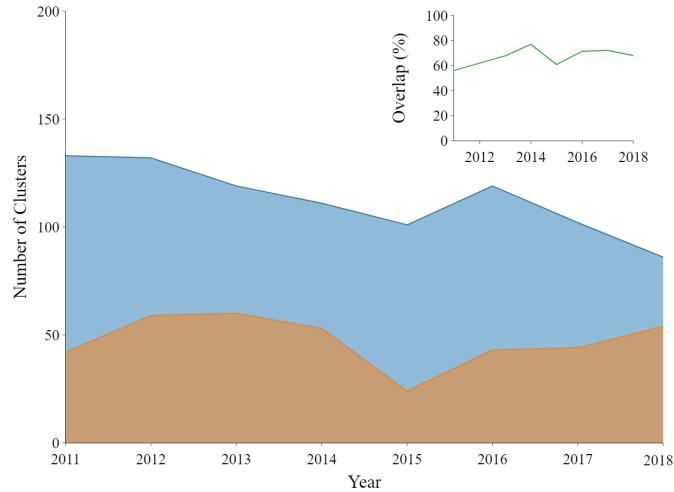
**Knowledge evolution is significantly associated with interdisciplinarity and weak members**

We use multinomial logistic regression with the mentioned endogenous and exogenous variables to evaluate how knowledge evolution may be explained through our interdisciplinarity scores, cluster size and network metrics. Per common practice, we insert a constant and a year variable to account for potential temporal effects. We attempt to explain whether or not clusters split or merge first, in order to evaluate the strength of associations between our hypothesized inputs and outputs.

Per Table 1, we see significant positive associations between a cluster splitting or merging and the language interdisciplinarity score and network interdisciplinarity score with only certain associations (i.e., without weak members).[6] We also see a positive association with the number of weak members associated with a cluster, and a negative association with the year.[7] Though all marginal effects are on the same order of magnitude, ranking by those effects, the language interdisciplinarity score is most important, followed by the number of weak members and the network score without weak members. Next, we further investigate this statistical relationship by testing the predictive power of a model trained on only a subset of these cluster data.

**Validating our statistical result with predictive power - equal importance of interdisciplinarity scores**
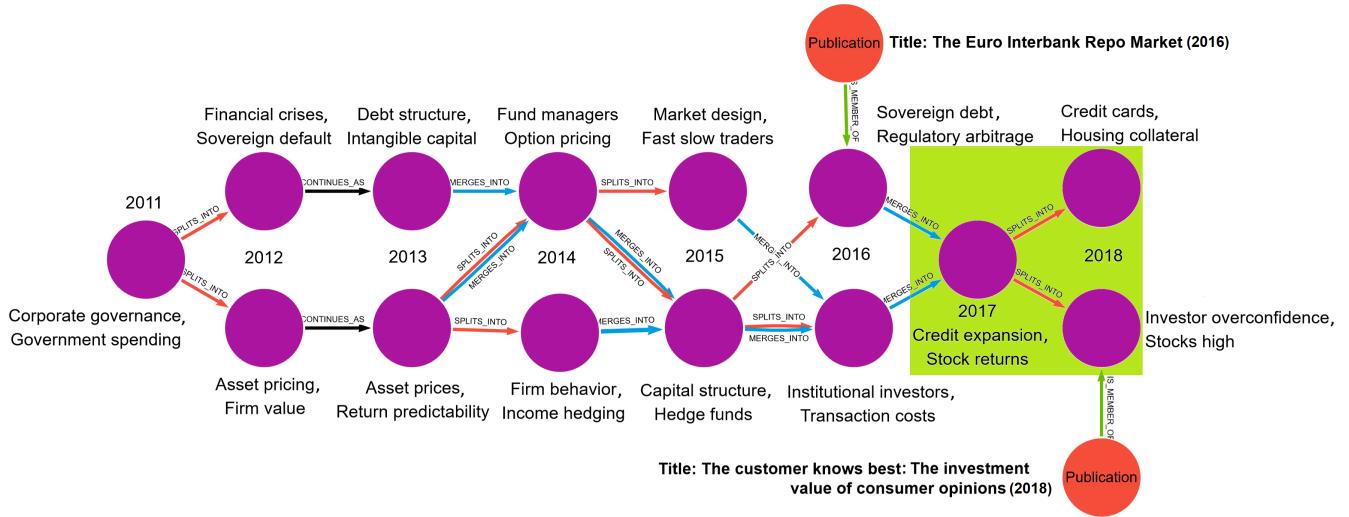
We have shown significant associations between knowledge splitting and merging, and interdisciplinarity and weak members. Here we go further by performing predictive modeling with a random forest classifier. Including only features shown to be statistically significant, we achieve a micro-averaged $F_1 = 0.67$ on our held-out test set, with $F_1 = 0.71$ on our class representing

---

[6]Following common best practice, we first conducted tests with all features, and, finding some insignificant, repeated with only significant features. See Supplementary Materials for details of this purposeful selection.

[7]Year was included per common practice to remove potential associations from time passing. Note that this model had a higher pseudo $R^2$ than a model without the year included. Future work should investigate any temporal associations through e.g., time series analyses.

**Figure 6.** Figure showing evolution of a set of language clusters from 2011 to 2018 (left to right) and keyphrases for each, along with two representative papers for two of the clusters. Note the marked change in focus between 2016 and 2018 evidenced by representative titles and cluster keyphrases. The split event for the 2017 cluster was successfully predicted by the random forest classifier described later (green box).

| | Model estimates | | Marginal effects | |
| --- | --- | --- | --- | --- |
| **Model Input (per cluster)** | **Coefficient** | *P* | **Effect** | *P* |
| Mean language ID score (strong members only) | 0.534 | 0.000 | 0.116 | 0.000 |
| Number of weak members | 0.449 | 0.003 | 0.097 | 0.002 |
| Mean network ID score (strong members only) | 0.292 | 0.030 | 0.063 | 0.025 |
| Publication year | -0.372 | 0.007 | -0.081 | 0.005 |
| Constant | -0.009 | 0.941 | | |

**Table 1.** Multinomial logistic regression results describing associations with split or merge (1) vs. continuation or death (0). Note significant positive associations with language score, network score, number of weak members, and a negative association for the year. All other features were not significant, and left out via purposeful selection for a more parsimonious model; see Supplemental Materials.

| Model input feature | Gini Importance |
| --- | --- |
| Mean language ID score (strong members only) | 0.336 |
| Number of weak members | 0.234 |
| Mean network ID score (strong members only) | 0.315 |
| Publication year | 0.115 |

**Table 2.** Random forest results on a held-out test set predicting the different types of cluster events a given cluster would experience in the next year, with the same features as in Table 1. We achieved a micro-averaged $F_1 = 0.67$ on our held-out test set, with a class-specific $F_1 = 0.71$ for the class representing knowledge evolution (splits and merges). Per reported Gini feature importance of each independent variable, both interdisciplinarity scores are equally important, followed by number of weak members, then year. Note that the sort order of this table is identical to that of Table 1 to allow for more direct comparison of logistic regression coefficients to random forest feature importances.

knowledge evolution (i.e. splitting or merging), a performance that is significantly better than random chance. Specifically, we present Table 2, which intuitively shows both interdisciplinarity scores to be equally important in achieving our predictive power. The number of weak members associated with a given cluster is next-most important, followed by the year variable. We validated against potential issues that can affect the Gini feature importance values from a random forest, specifically issues that arise when features exhibit multicollinearity and a bias towards numeric and high-cardinality categorical features[63]. The first is not a problem in this case, as the high-correlation features were removed as a result of the logistic regression analysis discussed earlier. The second is expected to only be a minimal concern for this analysis, as the only non-numeric feature in this model is the publication year. Because this is a low-cardinality categorical feature, it may be the victim of a bias in the feature importances and, as a result, the year's true ranking in the feature importance table could be higher than is indicated. As this is not a critical change in the data for our analysis, correcting for this bias is beyond the scope of this work. Taken together, our results underscore the importance of including both the linguistic and network viewpoints of interdisciplinarity.

## Conclusions and Future Work

In this paper, we proposed a hybrid language- and network-based framework that uses state-of-the-art semantic embeddings and citation information to model metastability of ideas in order to identify dynamic events associated with the rise, fall, combination, and dispersion of topics in the scholarly corpus. We show that this hybrid approach is distinctly different from those based on linguistic or citation information alone. The approach we propose relies on a multi-dimensional view of interdisciplinarity as a predictor of scientific knowledge transitions.

Through both explanatory and predictive efforts, we show that language as well as network interdisciplinarity has positive effects on metastable knowledge combining and mixing. Interestingly, network interdisciplinarity of strong member papers is significantly predictive of these mixing events, even though the number of strong members is not. By contrast, even though weak members' network interdisciplinarity is not significantly predictive, more weak members are predictive of knowledge combining and mixing. This suggests that papers that do not cluster neatly are indicative of combinatorial innovation that is expressed as the knowledge mixing events discussed herein. As such, if one is interested in spurring broad interdisciplinarity, one should focus on encouraging more weakly-clustered research, regardless of its own network-derived interdisciplinarity. Future work should further investigate these relationships, in particular over longer time scales and on a more complete body of the scientific corpus. Additionally, a comparison of a few other useful and popular interdisciplinarity metrics is a natural extension of this work, to determine how well other established measures can predict the knowledge evolution events we have explored. A lack of consensus on the most useful interdisciplinarity metrics[39] however makes this a challenge that must be tackled in later analyses.

This work also motivates and lays groundwork for new hybrid models that align multiple views of the literature ( e.g., linguistic, bibliometric) into unified modeling frameworks. Looking beyond traditional single-view approaches, such frameworks would be better suited to capture the richness of the scholarly record. This can be achieved through so-called graph machine learning modeling, that allows an integrated representation of a datum reflecting both its content (e.g. language in the case of a scientific paper) and its context within a network. Further, the work we describe here is mostly based on unsupervised learning. This is a necessity of the nature of this work, as there is no readily-available ground truth that is universally acknowledged to reflect the changing nature of scientific thought, disciplines, and sub-disciplines at a time scale reflective of how ideas mature and evolve. Future work should build benchmark datasets with which the metascience community can engage to evaluate and test these approaches more thoroughly than is currently possible. Possible proxy datasets that do exist at the moment include citation records – the prediction of above-average citation growth, for example, could be another modeling task that is able to further determine the utility of the interdisciplinarity metrics presented in this paper.

## Data Availability

Parts of the data that support the findings of this study are available from Clarivate but restrictions apply to the availability of these data, which were used under license for the current study, and so are not publicly available. Data are however available from the authors upon reasonable request and with permission of Clarivate.

## Code Availability

The code used for data processing and model development for the current study is available at https://github.com/QS-2/VESPID.git.

## References

1. Morris, C. The significance of the unity of science movement. *Philos. Phenomenol. Res.* **6**, 508–515 (1946).

2. Wang, D., Song, C. & Barabási, A.-L. Quantifying long-term scientific impact. *Science* **342**, 127–132 (2013).

3. Sinatra, R., Wang, D., Deville, P., Song, C. & Barabási, A.-L. Quantifying the evolution of individual scientific impact. *Science* **354** (2016).

4. Liu, L. *et al.* Hot streaks in artistic, cultural, and scientific careers. *Nature* **559**, 396–399 (2018).

5. Li, J., Yin, Y., Fortunato, S. & Wang, D. Scientific elite revisited: patterns of productivity, collaboration, authorship and impact. *J. Royal Soc. Interface* **17**, 20200135 (2020).

6. Pluchino, A. *et al.* Exploring the role of interdisciplinarity in physics: success, talent and luck. *PloS one* **14**, e0218793 (2019).

7. Janosov, M., Battiston, F. & Sinatra, R. Success and luck in creative careers. *EPJ Data Sci.* **9**, 1–12 (2020).

8. Peterson, D. & Panofsky, A. Metascience as a scientific social movement. *SocArXiv* (2020).

9. Schooler, J. W. Metascience could rescue the 'replication crisis'. *Nat. News* **515**, 9 (2014).

10. Larivière, V., Ni, C., Gingras, Y., Cronin, B. & Sugimoto, C. R. Bibliometrics: Global gender disparities in science. *Nat. News* **504**, 211 (2013).

11. Hofstra, B. *et al.* The diversity–innovation paradox in science. *Proc. Natl. Acad. Sci.* **117**, 9284–9291 (2020).

12. Franco, A., Malhotra, N. & Simonovits, G. Publication bias in the social sciences: Unlocking the file drawer. *Science* **345**, 1502–1505 (2014).

13. Rzhetsky, A., Foster, J. G., Foster, I. T. & Evans, J. A. Choosing experiments to accelerate collective discovery. *Proc. Natl. Acad. Sci.* **112**, 14569–14574 (2015).

14. Jia, T., Wang, D. & Szymanski, B. K. Quantifying patterns of research-interest evolution. *Nat. Hum. Behav.* **1**, 1–7 (2017).

15. Spangler, S. *et al.* Automated hypothesis generation based on mining scientific literature. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 1877–1886 (2014).

16. Prabhakaran, V., Hamilton, W. L., McFarland, D. & Jurafsky, D. Predicting the rise and fall of scientific topics from trends in their rhetorical framing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1170–1180 (2016).

17. Chen, C., Wang, Z., Li, W. & Sun, X. Modeling scientific influence for research trending topic prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32 (2018).

18. Kleminski, R. & Kazienko, P. Identifying promising research topics in computer science. In *European Network Intelligence Conference*, 231–241 (Springer, 2017).

19. Uban, A. S., Caragea, C. & Dinu, L. P. Studying the evolution of scientific topics and their relationships. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 1908–1922 (2021).

20. Faust, O. Documenting and predicting topic changes in computers in biology and medicine: A bibliometric keyword analysis from 1990 to 2017. *Informatics Medicine Unlocked* **11**, 15–27 (2018).

21. Shibata, N., Kajikawa, Y., Takeda, Y. & Matsushima, K. Detecting emerging research fronts based on topological measures in citation networks of scientific publications. *Technovation* **28**, 758–775 (2008).

22. Salatino, A. A., Osborne, F. & Motta, E. Augur: forecasting the emergence of new research topics. In *Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries*, 303–312 (2018).

23. Sun, Y. & Latora, V. The evolution of knowledge within and across fields in modern physics. *Sci. reports* **10**, 1–9 (2020).

24. Zamani, M. *et al.* Evolution and transformation of early modern cosmological knowledge: A network study. *Sci. Reports* **10**, 1–15 (2020).

25. Sarigöl, E., Pfitzner, R., Scholtes, I., Garas, A. & Schweitzer, F. Predicting scientific success based on coauthorship networks. *EPJ Data Sci.* **3**, 1–16 (2014).

26. Sun, X., Ding, K. & Lin, Y. Mapping the evolution of scientific fields based on cross-field authors. *J. Informetrics* **10**, 750–761 (2016).

27. Krenn, M. & Zeilinger, A. Predicting research trends with semantic and neural networks with an application in quantum physics. *Proc. Natl. Acad. Sci.* **117**, 1910–1916 (2020).

28. Sasaki, H., Fugetsu, B. & Sakata, I. Emerging scientific field detection using citation networks and topic models—a case study of the nanocarbon field. *Appl. Syst. Innov.* **3**, 40 (2020).

29. Zhang, Y., Chen, H., Lu, J. & Zhang, G. Detecting and predicting the topic change of knowledge-based systems: A topic-based bibliometric analysis from 1991 to 2016. *Knowledge-Based Syst.* **133**, 255–268 (2017).

30. Cohan, A., Feldman, S., Beltagy, I., Downey, D. & Weld, D. S. Specter: Document-level representation learning using citation-informed transformers. *arXiv preprint arXiv:2004.07180* (2020).

31. Klein, J. T. *Interdisciplinarity: History, theory, and practice* (Wayne state university press, 1990).

32. Jacobs, J. A. & Frickel, S. Interdisciplinarity: A critical assessment. *Annu. review Sociol.* **35**, 43–65 (2009).

33. Repko, A. F., Szostak, R. & Buchberger, M. P. *Introduction to interdisciplinary studies* (Sage Publications, 2019).

34. Pan, R. K., Sinha, S., Kaski, K. & Saramäki, J. The evolution of interdisciplinarity in physics research. *Sci. reports* **2**, 1–8 (2012).

35. Molas-Gallart, J., Rafols, I. & Tang, P. On the relationship between interdisciplinarity and impact: Different modalities of interdisciplinarity lead to different types of impact (< special report> toward interdisciplinarity in research and development). *The J. Sci. Policy Res. Manag.* **29**, 69–89 (2014).

36. Jacobs, J. A. *In defense of disciplines* (University of chicago Press, 2014).

37. Wagner, C. S. *et al.* Approaches to understanding and measuring interdisciplinary scientific research (idr): A review of the literature. *J. informetrics* **5**, 14–26 (2011).

38. Porter, A. & Rafols, I. Is science becoming more interdisciplinary? measuring and mapping six research fields over time. *Scientometrics* **81**, 719–745 (2009).

39. Wang, Q. & Schneider, J. W. Consistency and validity of interdisciplinarity measures. *Quant. Sci. Stud.* **1**, 239–263 (2020).

40. Analytics, C. Web of science (2021).

41. Porter, A. & Chubin, D. An indicator of cross-disciplinary research. *Scientometrics* **8**, 161–176 (1985).

42. Morillo, F., Bordons, M. & Gómez, I. An approach to interdisciplinarity through bibliometric indicators. *Scientometrics* **51**, 203–222 (2001).

43. Wang, J., Thijs, B. & Glänzel, W. Interdisciplinarity and impact: Distinct effects of variety, balance, and disparity. *PloS one* **10**, e0127298 (2015).

44. Leydesdorff, L. Betweenness centrality as an indicator of the interdisciplinarity of scientific journals. *J. Am. Soc. for Inf. Sci. Technol.* **58**, 1303–1319 (2007).

45. Leydesdorff, L. & Rafols, I. Indicators of the interdisciplinarity of journals: Diversity, centrality, and citations. *J. Informetrics* **5**, 87–100 (2011).

46. Gao, Q. *et al.* Potential index: Revealing the future impact of research topics based on current knowledge networks. *J. Informetrics* **15**, 101165 (2021).

47. Chen, C. *et al.* Towards an explanatory and computational theory of scientific discovery. *J. Informetrics* **3**, 191–209 (2009).

48. Alipourfard, N. *et al.* Systematizing confidence in open research and evidence (score). *SocArXiv* (2021).

49. Ammar, W. *et al.* Construction of the literature graph in semantic scholar. *arXiv preprint arXiv:1805.02262* (2018).

50. Lo, K., Wang, L. L., Neumann, M., Kinney, R. & Weld, D. S. S2orc: The semantic scholar open research corpus. *arXiv preprint arXiv:1911.02782* (2019).

51. Lammey, R. Crossref text and data mining services. *Insights* **28** (2015).

52. Wolf, T. *et al.* Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45 (Association for Computational Linguistics, Online, 2020).

53. McInnes, L., Healy, J. & Astels, S. hdbscan: Hierarchical density based clustering. *J. Open Source Softw.* **2**, 205 (2017).

54. McInnes, L., Healy, J. & Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).

55. Turner, R. *et al.* Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In *NeurIPS 2020 Competition and Demonstration Track*, 3–26 (PMLR, 2021).

56. Grootendorst, M. Keybert: Minimal keyword extraction with bert., DOI: 10.5281/zenodo.4461265 (2020).

57. Lu, H., Halappanavar, M. & Kalyanaraman, A. Parallel heuristics for scalable community detection. *Parallel Comput.* **47**, 19–37 (2015).

58. Lancichinetti, A. & Fortunato, S. Community detection algorithms: a comparative analysis. *Phys. review E* **80**, 056117 (2009).

59. Rafols, I. & Meyer, M. Diversity measures and network centralities as indicators of interdisciplinarity: case studies in bionanoscience. In *Proceedings of ISSI*, vol. 2, 631–637 (2007).

60. Greene, D., Doyle, D. & Cunningham, P. Tracking the evolution of communities in dynamic social networks. In *2010 international conference on advances in social networks analysis and mining*, 176–183 (IEEE, 2010).

61. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).

62. Asur, S., Parthasarathy, S. & Ucar, D. An event-based framework for characterizing the evolutionary behavior of interaction graphs. *ACM Transactions on Knowl. Discov. from Data (TKDD)* **3**, 1–36 (2009).

63. Strobl, C., Boulesteix, A.-L., Zeileis, A. & Hothorn, T. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics* **8**, 1–21 (2007).

## Acknowledgements

## Author contributions statement

D.G., S.R. conceived the experiment(s); S.K., D.R.M. and M.S. conducted the experiment(s). All authors analyzed the results and reviewed the manuscript.

## Additional information

The authors declare no competing interests.