



HAL
open science

A Novel Interdisciplinarity Model Towards Inter-domain Information Pairing

Nicolas Douard, Ahmed Samet, George Giakos, Denis Cavallucci

► **To cite this version:**

Nicolas Douard, Ahmed Samet, George Giakos, Denis Cavallucci. A Novel Interdisciplinarity Model Towards Inter-domain Information Pairing. TFC 2024, Cluj-Napoca, Romania, novembre 2024, Nov 2024, Cluj-Napoca, Romania. hal-04922408

HAL Id: hal-04922408

<https://hal.science/hal-04922408v1>

Submitted on 30 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Novel Interdisciplinarity Model Towards Inter-Domain Information Pairing

Nicolas Douard^{1,2}, Ahmed Samet¹, George Giakos², and Denis Cavallucci¹

¹ National Institute of Applied Sciences, University of Strasbourg,
24 Bd de la Victoire, 67000 Strasbourg, France

² Department of Electrical and Computer Engineering, Manhattan College,
3825 Corlear Ave, New York, NY 10463, USA

Abstract. This study introduces an interdisciplinary prediction framework as part of a novel approach that integrates the Inventive Design Method (IDM), Topic Modeling, and Generative AI to foster innovation across academic fields. Identifying interdisciplinary connections is essential for solving complex, multi-domain problems. Our research uses a supervised machine learning classifier to identify interdisciplinary documents within the Semantic Scholar corpus, extracting latent insights. The Text Convolutional Neural Network model performed best, achieving an F1 score of 0.80. We find that approximately 25% of human knowledge is interdisciplinary. This framework helps create comprehensive knowledge maps across multiple domains, promoting innovation through effective cross-domain knowledge transfer.

Keywords: Interdisciplinary mapping · cross-domain transfer · biomimicry · literature analysis · topic models · language models.

1 Introduction

Biomimicry leverages nature's strategies to design and produce materials, structures, and systems that solve human challenges sustainably. This research aims to unlock the untapped potential of the natural world, providing novel and sustainable solutions to complex engineering problems.

Building upon the groundwork presented in "TFC 2023" [1], this study introduces a novel framework that combines the IDM methodology [2] with Generative AI. This innovative approach enhances the retrieval and application of bioinspired solutions for engineering problems [3], enabling targeted extraction of relevant articles from extensive databases [4].

Interdisciplinarity integrates concepts, theories, and methodologies from different disciplines to solve complex problems, create new knowledge, and achieve a comprehensive understanding. Cross-domain knowledge transfer applies insights, methods, or technologies from one field to another, fostering innovation and solving problems that span multiple expertise areas.

Identifying interdisciplinary connections across diverse knowledge domains is crucial for driving significant innovation, yet the vast scale of scientific literature

renders manual mapping impractical. Traditional methods, such as bibliometrics (e.g., co-authorships, collaborations, citations) and network dynamics (e.g., betweenness centrality, diversity) [5], while valuable, demand sophisticated interpretations and qualitative assessments that lack scalability. In response, we introduce a computationally scalable method aiming to systematically extract latent interdisciplinary insights from the comprehensive corpus of human knowledge.

Inspired by Altshuller’s work on TRIZ [6] and derivatives [2], knowledge can be structured as follows:

- Case C1: the solution is within the same industry
- Case C2: the solution is in another industry
- Case C3: the solution outside of what exists in all industries
- Case C4: the solution does not yet exist

A schematic representation of these cases is shown in Fig. 1, with this research focusing on case C2.

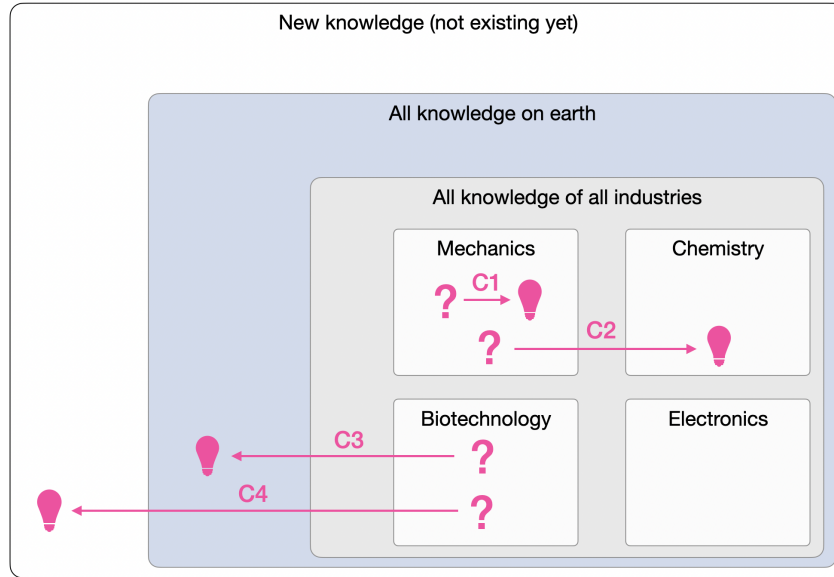


Fig. 1. Frames of knowledge.

A comprehensive dataset of scientific literature is essential to uncover interdisciplinary associations within existing knowledge. Preliminary investigations revealed that simple keyword mapping, as shown in Fig. 2, inadequately captures latent thematic structures in textual data [7]. Advanced unsupervised semantic

modeling techniques are needed to fully reveal these complex inter-domain connections [8] [1].

Elucidating interdisciplinary associations within scientific knowledge requires a substantial, discipline-spanning corpus. However, undifferentiated monodisciplinary and multidisciplinary texts can lead to modeling challenges, such as overfitting on irrelevant features and obscuring key thematic structures with noise. Differentiating corpora is crucial to obtain an optimized training set that highlights multidisciplinary connections and mitigates the risks of producing diluted, inaccurate models.

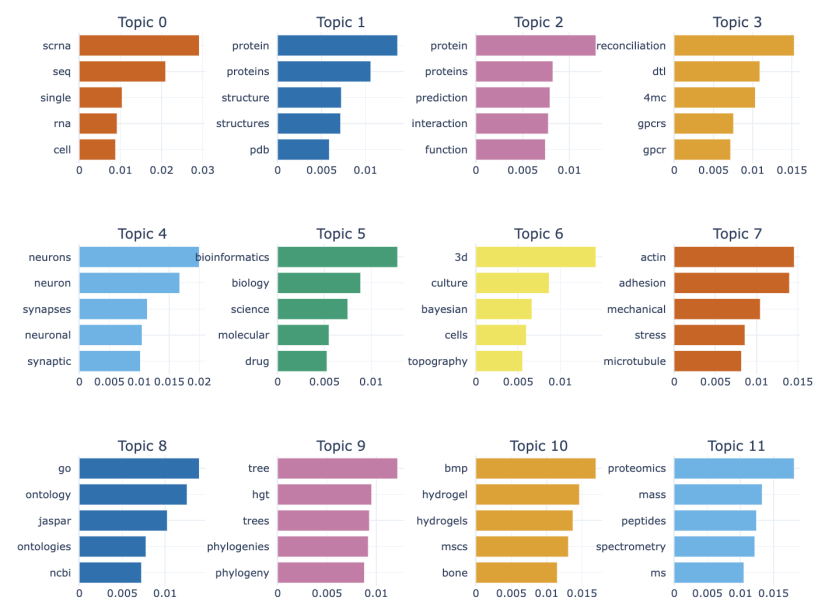


Fig. 2. Topic word scores snippet.

2 Experiment

2.1 Materials and Methods

To identify texts that span multiple disciplines rather than being confined to a single field, we propose using a supervised machine learning classifier. This classifier aims at recognizing interdisciplinary documents, even when metadata labels them under just one domain. Accurately identifying such interdisciplinary works is necessary for constructing comprehensive maps of knowledge connections across disciplines.

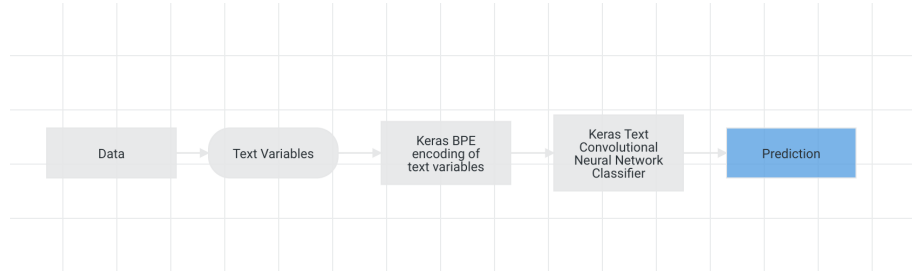


Fig. 3. Overview of the proposed model architecture.

Unlike traditional bibliometric approaches that rely heavily on explicit meta-data and may overlook interdisciplinary works [5], our method leverages full-text content and advanced semantic modeling to detect latent interdisciplinary connections, thereby addressing the scalability issues inherent in existing methods.

For developing interdisciplinary knowledge maps, this classifier serves as an initial filter to extract relevant interdisciplinary documents for subsequent topic modeling and other analyses. Filtering out purely discipline-specific documents reduces noise, allowing the downstream methods to focus on substantively interdisciplinary content.

The overall methodology prioritizes concentrating on key information that bridges disciplines, guided by TRIZ [6] principles of resolving contradictions through an interdisciplinary lens.

2.2 Dataset

Semantic Scholar’s corpus of over 100,000,000 abstracts [4] encompasses a very large breadth of disciplines, providing a substantial foundation for analysis.

While the Semantic Scholar corpus metadata labels each document’s primary discipline(s), this categorization is imperfect - many articles classified in one field may still relate to other domains. Metadata labels alone are insufficient for comprehensively identifying interdisciplinary works.

The combination of full-text content and associated metadata enables the development of predictive models that can identify interdisciplinary connections beyond the primary discipline labels. This approach is based on the fundamental assumption that such models can effectively generalize patterns across disciplines.

Semantic Scholar data is downloaded using the bulk API [4] and pre-processed to create a target column indicating whether a given article is marked as belonging to one or several *fields of study*.

There are 18 referenced fields of study in the Semantic Scholar database. Given focus on engineering, physics, and TRIZ [6] scope, the following can be deemed of interest:

- Medicine
- Computer Science
- Chemistry
- Physics
- Mathematics
- Engineering
- Biology
- Materials Science

2.3 Model

Different estimators are evaluated:

- Text Convolutional Neural Network [9]
- Gradient Boosted Trees Classifier with Early Stopping (Fast Feature Binning) [10] [11]
- Elastic-Net Classifier [12]
- Deep Residual Neural Network Classifier using Training Schedule (3 Layers: 512, 64, 64 Units) [13]

Learning curves are plotted to evaluate model performance, in terms of cross-entropy loss (or LogLoss), as the sample size changes, as shown in Fig. 4.

Interestingly enough, Text Convolutional Neural Networks (Text CNNs) use convolutional layers to process textual data, capturing local features. By employing filters of various sizes, Text CNNs can detect patterns at multiple scales, enhancing semantic understanding. Additionally, the use of pooling layers reduces dimensionality, and dropout layers prevent overfitting, making these models robust and generalizable across different text datasets [9].

Gradient boosting works by combining several weak predictive models into a single strong model. In the context of this architecture, gradient boosting is used to combine different features into a single model [10] that can be used for classification. Early stopping in gradient boosting helps identify the minimum iterations needed to develop a model that generalizes well [11] while preventing overfitting. This is achieved by tracking validation accuracy after each iteration and stopping once little improvement is seen. The model becomes simpler yet more robust by reducing overfitting and complexity.

The Elastic Net classifier [12] combines the strengths of both Ridge and Lasso regression methods to enhance model performance, especially when dealing with highly correlated data. This hybrid model employs a mix of L1 and L2 regularization to penalize complex models, encouraging sparsity while also maintaining model robustness. The balance between L1 and L2 penalties can be adjusted to suit specific datasets, making Elastic Net highly adaptable. This method is effective in preventing overfitting and is useful for variable selection in scenarios with numerous predictors.

The Deep Residual Neural Network (ResNet) classifier [13], utilizing a training schedule with three layers of 512, 64, and 64 units, addresses the vanishing

gradient problem in deep neural networks. This architecture incorporates residual blocks that allow activations to skip one or more layers, promoting smoother training and convergence. The specified training schedule and layer configuration enable the network to learn complex patterns while maintaining a manageable computational load.

3 Results

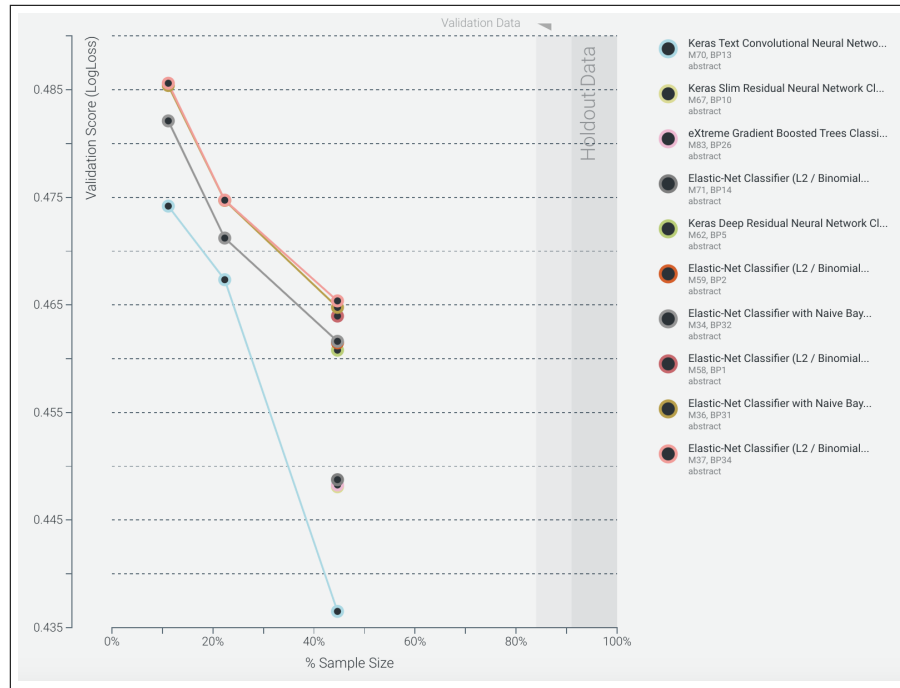


Fig. 4. Learning curves for key models inform on performance variation as the sample size changes.

The 1D Text Convolutional Neural Network approach produced the lowest cross-entropy loss score, outperforming other contenders. Specifically, the model exhibits a true positive rate (sensitivity) of 0.84 and a positive predictive value (precision) of 0.76 which yields an F1 score of 0.80. A schematic representation of this pipeline is shown in Fig. 3.

The model prediction distribution is shown in Fig. 5. On the graph, red represents values that fall into the unsuitable or negative class, while blue represents values that fall into the suitable or positive class. The prediction threshold was set to maximize Matthew’s Correlation Coefficient (MCC) [14].

Inference over the full Semantic Scholar dataset suggests that about 25% of all human knowledge can be seen as interdisciplinary.

The prevalence of interdisciplinary knowledge uncovered by this analysis underscores both the challenges and immense potential in systematically mapping these intricate inter-domain links. The diffuse, decentralized nature of interdisciplinary insights makes their manual consolidation impractical. However, automated methods that can comprehensively survey and distill interdisciplinary patterns from massive datasets offer a compelling solution.

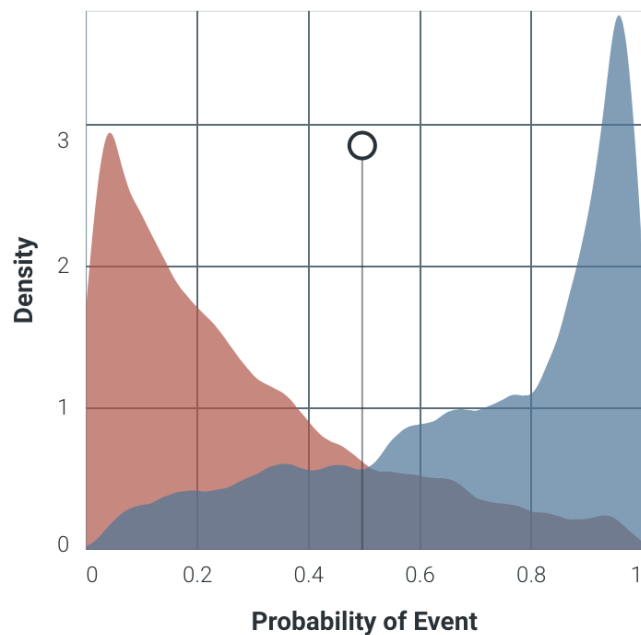


Fig. 5. Distribution of predictions. Prediction threshold tuned to maximize Matthew's Correlation Coefficient.

The Text Convolutional Neural Network has shown promise in identifying interdisciplinary documents. However, this approach has notable limitations. Its dependence on text data and metadata may introduce biases, particularly in language and disciplinary terminology. This could result in overlooking subtle interdisciplinary connections. Moreover, the model's ability to maintain accuracy when scaled up remains an important assumption that requires further validation.

4 Conclusion

This research presents a machine learning approach to automatically identify interdisciplinary documents from a large academic corpus. By employing a text convolutional neural network classifier, the methodology distinguishes interdisciplinary texts, even when metadata classifies them within a single domain. Precisely pinpointing interdisciplinary documents is a crucial first step towards constructing comprehensive interdisciplinary knowledge maps and enabling systematic cross-domain transfer. The ability to isolate an interdisciplinary subset of literature allows subsequent topic modeling and language AI techniques to be precisely targeted, facilitating the discovery of latent links across diverse fields.

This interdisciplinary document prediction framework lays the foundation for mapping connections between domains, which has immense potential for driving innovation through biomimicry and interdisciplinary insights, as explored in "TFC 2022" [7]. Future research can build upon this critical filtering step to develop fuller interdisciplinary mapping and knowledge transfer pipelines. The integration of this proposed framework into existing academic tools and platforms can also be explored alongside hybrid approaches where AI-driven models are complemented by expert validation.

Bibliography

- [1] Nicolas Douard, Ahmed Samet, George C. Giakos, and Denis Cavallucci. Navigating the knowledge network: How inter-domain information pairing and generative ai can enable rapid problem-solving. In *TFC*, 2023.
- [2] Denis Cavallucci. From triz to inventive design method (idm): towards a formalization of inventive practices in r&d departments. 2012.
- [3] Meredith Ringel Morris. Scientists’ perspectives on the potential for generative ai in their fields. *ArXiv*, abs/2304.01420, 2023.
- [4] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Michael Kinney, and Daniel S. Weld. S2orc: The semantic scholar open research corpus. In *Annual Meeting of the Association for Computational Linguistics*, 2020.
- [5] Caroline S. Wagner, J. David Roessner, Kamau Bobb, Julie Thompson Klein, Kevin W. Boyack, Joann Keyton, Ismael Ràfols, and Katy Börner. Approaches to understanding and measuring interdisciplinary scientific research (idr): A review of the literature. *J. Informetrics*, 5:14–26, 2011.
- [6] G. S. Altshuller and Lev A. Shulyak. *And Suddenly the Inventor Appeared: TRIZ, the Theory of Inventive Problem Solving*. Technical Innovation Center, 1996.
- [7] Nicolas Douard, Ahmed Samet, George C. Giakos, and Denis Cavallucci. Bridging two different domains to pair their inherent problem-solution text contents: Applications to quantum sensing and biology. In *TFC*, 2022.
- [8] Maarten R. Grootendorst. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *ArXiv*, abs/2203.05794, 2022.
- [9] Yajian Zhou, Jiale Li, Junhui Chi, Wei Tang, and Yuqi Zheng. Set-cnn: A text convolutional neural network based on semantic extension for short text classification. *Knowledge-Based Systems*, 257:109948, 2022.
- [10] Alexey Natekin and Alois Knoll. Gradient boosting machines, a tutorial. *Frontiers in Neurorobotics*, 7, 2013.
- [11] Andreas Mayr, Benjamin Hofner, and Matthias Schmid. The importance of knowing when to stop. *Methods of Information in Medicine*, 51:178–186, 2012.
- [12] Hui Zou and Trevor J. Hastie. Addendum: Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67, 2005.
- [13] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2015.
- [14] Davide Chicco and Giuseppe Jurman. The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. *BMC Genomics*, 21, 2020.