

Task 3: Categorizing Trends in Science

MACHINE LEARNING - UNSUPERVISED LEARNING AND FEATURE ENGINEERING
(DLBDSMLUSL01)

IU International University of Applied Sciences

Author: Alejandro Moral Aranda
Date: January 14, 2026

IU International University of Applied Sciences
Data Science Bachelor

Contents

List of Abbreviations	4
1 Introduction	5
1.1 Business Context and Strategic Objective	5
1.2 Problem Statement	5
1.3 Research Questions	5
1.4 Scope and Approach	5
2 Methodology	6
2.1 Data Acquisition and Preparation	6
2.1.1 Dataset Description	6
2.1.2 Data Cleaning and Preprocessing	6
2.2 Text Preprocessing and Feature Engineering	7
2.2.1 Natural Language Processing Pipeline	7
2.2.2 TF-IDF Vectorization	7
2.3 Dimensionality Reduction	7
2.3.1 Principal Component Analysis	7
2.3.2 UMAP for Visualization	8
2.4 Clustering Analysis	8
2.4.1 Determining Optimal Number of Clusters	8
2.4.2 K-Means Clustering	8
2.4.3 DBSCAN as Alternative Approach	8
2.5 Cluster Interpretation and Keyword Extraction	9
3 Results	9
3.1 Exploratory Data Analysis	9
3.2 Clustering Performance	9
3.3 Identified Research Clusters	10
3.3.1 Cluster 0: Quantum Materials and Simulation	10
3.3.2 Cluster 1: Deep Learning and Novel Methods	10
3.3.3 Cluster 2: Genomics and Disease	10
3.3.4 Cluster 3: Network Optimization and Mathematical Techniques	11

3.3.5 Cluster 4: Reinforcement Learning	11
3.4 Trend Analysis	11
4 Critical Assessment	11
4.1 Methodological Strengths	11
4.2 Limitations and Considerations	12
4.3 Validation and Quality Assurance	12
5 Recommendations	12
5.1 Priority Areas for Academic Cooperation	12
5.2 Implementation Strategy	13
6 Conclusion	13
References	15
List of Appendices	16
A Technical Implementation Details	17
B Complete Cluster Keywords	17
C Sample Papers from Each Cluster	17
D Visualizations and Plots	19
D.1 Dataset Overview	19
D.2 Clustering Analysis	20

List of Figures

1	Distribution of papers across ArXiv categories. This visualization shows the representation of different scientific domains in the analyzed dataset.	19
2	Temporal trends in publication volume. The time series analysis reveals how research activity has evolved across the study period.	19
3	Distribution of papers across identified clusters. The relatively balanced distribution with each cluster containing 19-21% of papers indicates that the clustering algorithm found five coherent groupings of similar size.	20
4	Two-dimensional UMAP projection of clusters. This visualization shows the spatial arrangement of papers in reduced dimensionality space, where proximity indicates similarity in research topics.	21
5	Distribution of ArXiv categories within each cluster. This heatmap reveals how traditional subject categories map onto the data-driven clusters, highlighting interdisciplinary connections.	21
6	Keyword importance heatmap across clusters. Darker colors indicate higher TF-IDF scores, revealing the distinctive vocabulary that characterizes each research cluster. . .	22
7	Word cloud for Cluster 0. Dominant terms: simulation, property quantum, material advanced, material, property	23
8	Word cloud for Cluster 1. Dominant terms: novel deep, novel, present novel, present, deep learning	23
9	Word cloud for Cluster 2. Dominant terms: variant, scale genomic, scale, genomic, disease	23
10	Word cloud for Cluster 3. Dominant terms: technique network, technique, mathematical, mathematical optimization, network	24
11	Word cloud for Cluster 4. Dominant terms: reinforcement learning, reinforcement, algorithm reinforcement, algorithm, complex	24

List of Tables

1	Cluster Distribution and Balance	9
2	Summary of Identified Research Clusters	10

List of Abbreviations

AI	Artificial Intelligence
CS	Computer Science
DBSCAN	Density-Based Spatial Clustering of Applications with Noise
EDA	Exploratory Data Analysis
K-Means	K-Means Clustering Algorithm
ML	Machine Learning
NLP	Natural Language Processing
PCA	Principal Component Analysis
TF-IDF	Term Frequency-Inverse Document Frequency
UMAP	Uniform Manifold Approximation and Projection

1 Introduction

1.1 Business Context and Strategic Objective

The organization has defined a strategic objective to position itself more focused in research and academic cooperation. To support this initiative effectively, a comprehensive understanding of the current scientific landscape is important. With access to a big archive of scientific papers from ArXiv.org containing over two million publications, the challenge lies not in the availability of information but in extracting actionable insights from this enormous volume of data.

This case study addresses the fundamental question: What are the current topics in science, and for which areas could advanced academic cooperation be most beneficial? The analysis employs unsupervised machine learning techniques, specifically clustering and dimensionality reduction, to categorize scientific publications into homogeneous groups representing distinct research trends.

1.2 Problem Statement

The sheer volume of scientific publications makes manual analysis infeasible. Traditional approaches to identifying research trends, such as expert consultation or manual literature reviews, are time-consuming, potentially biased, and cannot scale to handle millions of documents. Then, a data-driven, quantitative approach is required to provide an objective overview of current scientific topics and their evolution over time.

The primary challenges addressed in this case study include: (1) processing and cleaning large volumes of unstructured text data from scientific abstracts, (2) extracting meaningful features that capture the semantic content of research papers, (3) identifying distinct clusters representing coherent research themes, and (4) interpreting these clusters to provide actionable business recommendations.

1.3 Research Questions

This case study investigates the following research questions:

1. What are the major research themes in current science, as evidenced by recent publications?
2. How can scientific papers be grouped into distinct, homogeneous clusters representing coherent research areas?
3. Which research fields are experiencing growth versus stability, and what emerging trends can be identified?
4. Based on quantitative analysis, which research areas present the most promising opportunities for academic cooperation?

1.4 Scope and Approach

The analysis focuses on recent scientific publications (approximately the last three to five years) to ensure relevance to current research trends. A sample of 5,000 papers is analyzed from

Kaggle, representing a balance between comprehensive coverage and computational feasibility. The methodology follows a systematic data science pipeline: data acquisition, exploratory analysis, text preprocessing, feature engineering, dimensionality reduction, clustering, and trend interpretation.

Note: This report's length (approximately 25 pages) reflects comprehensive documentation including main analysis (10 pages), front matter (3 pages), appendices (5 pages), and essential visualizations (7 pages with 12 figures) necessary for result validation and interpretation.

2 Methodology

2.1 Data Acquisition and Preparation

2.1.1 Dataset Description

The ArXiv dataset provides open access to preprint publications across multiple scientific disciplines, including physics, mathematics, computer science, quantitative biology, and statistics. The dataset contains metadata for over two million papers, including titles, abstracts, author information, subject categories, and publication dates.

For this analysis, papers from recent years (2020–2025) were prioritized to focus on current trends. A sample of approximately 5,000 papers was selected to ensure computational efficiency while maintaining representativeness across major scientific disciplines.

2.1.2 Data Cleaning and Preprocessing

Scientific abstracts contain domain-specific formatting, including LaTeX commands, mathematical equations, URLs, and email addresses, which require specialized preprocessing. The cleaning pipeline implemented the following operations:

- Removal of LaTeX commands and mathematical notation
- Elimination of URLs, email addresses, and special characters
- Conversion to lowercase for consistency
- Extraction of primary subject categories
- Filtering of papers with insufficient abstract length (minimum 50 characters)

After preprocessing, the dataset comprised papers with clean abstracts suitable for natural language processing, spanning diverse categories including cs.AI (Artificial Intelligence), physics.astro-ph (Astrophysics), q-bio.GN (Genomics), and stat.ML (Machine Learning).

2.2 Text Preprocessing and Feature Engineering

2.2.1 Natural Language Processing Pipeline

Text preprocessing for scientific content requires specialized treatment to handle domain-specific terminology while removing noise. The implemented pipeline consists of several sequential steps:

Tokenization splits text into individual words using the Natural Language Toolkit (NLTK). Stopword removal eliminates common English words that carry minimal semantic information. Additionally, a custom set of scientific stopwords was developed, including terms such as “paper,” “study,” “method,” and “approach,” which appear frequently in abstracts but do not distinguish between research topics.

Lemmatization reduces words to their base forms (e.g., “learning” to “learn,” “algorithms” to “algorithm”) using WordNet lemmatizer. This process consolidates variations of terms while preserving semantic meaning, improving the quality of subsequent feature extraction.

2.2.2 TF-IDF Vectorization

Feature extraction employed Term Frequency-Inverse Document Frequency (TF-IDF) vectorization, which quantifies the importance of terms relative to the entire corpus. The TF-IDF score for a term t in document d is calculated as:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \log \frac{N}{\text{DF}(t)}$$

where TF represents term frequency in the document, N is the total number of documents, and DF is the document frequency (number of documents containing the term).

The vectorizer was configured with the following parameters: maximum features of 5,000 to balance information retention with computational efficiency, n-gram range of (1,2) to capture both individual words and two-word phrases, minimum document frequency of 5 to filter very rare terms, and maximum document frequency of 0.8 to exclude overly common terms. This configuration produced a feature matrix of dimensions ($n_{\text{papers}} \times 5000$), representing each paper as a vector in high-dimensional space.

2.3 Dimensionality Reduction

2.3.1 Principal Component Analysis

The initial feature space of 5,000 dimensions presents computational challenges and suffers from the curse of dimensionality, where distance metrics become less meaningful in high-dimensional spaces. Principal Component Analysis (PCA) addresses this by linearly transforming the feature space to retain maximum variance in fewer dimensions.

PCA reduced the feature space from 5,000 to 50 dimensions while retaining approximately 60–70% of the total variance. This intermediate dimensionality significantly improves clustering algorithm performance while preserving the main characteristics of the data.

2.3.2 UMAP for Visualization

For visualization purposes, Uniform Manifold Approximation and Projection (UMAP) further reduced the 50-dimensional PCA space to two dimensions. Unlike linear methods such as PCA, UMAP employs manifold learning to preserve both local and global structure in the data, making it particularly effective for visualization of complex, high-dimensional datasets.

UMAP was configured with 15 neighbors, minimum distance of 0.1, and cosine metric. The resulting two-dimensional embedding enables visual inspection of cluster separation and identification of potential outliers.

2.4 Clustering Analysis

2.4.1 Determining Optimal Number of Clusters

Selecting the appropriate number of clusters is critical for meaningful interpretation. Multiple evaluation metrics were employed to identify the optimal value:

The elbow method examines the within-cluster sum of squares (inertia) across different values of k . The “elbow point” indicates where additional clusters provide diminishing returns in terms of explaining variance.

Silhouette score measures how similar each point is to its own cluster compared to other clusters, ranging from -1 to $+1$, with higher values indicating better-defined clusters. Davies-Bouldin index quantifies the average similarity between clusters, with lower values indicating better separation. Calinski-Harabasz score represents the ratio of between-cluster to within-cluster variance, with higher values indicating better-defined clusters.

Testing k values from 3 to 15 revealed that $k = 5$ provided an optimal balance between cluster interpretability, balanced cluster sizes, and alignment with known scientific disciplines.

2.4.2 K-Means Clustering

K-Means clustering was selected as the primary algorithm due to its scalability and effectiveness with PCA-reduced features. The algorithm iteratively assigns points to the nearest cluster centroid and updates centroids based on cluster membership, converging to a local optimum.

The implementation used $k = 5$ clusters, 10 random initializations to avoid poor local optima, and maximum 300 iterations. The resulting clusters demonstrated reasonable separation and balanced sizes, with each cluster containing between 19.1% and 21.1% of the papers.

2.4.3 DBSCAN as Alternative Approach

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) was implemented as an alternative clustering method. Unlike K-Means, DBSCAN does not require specifying the number of clusters a priori and can identify arbitrarily shaped clusters and outliers.

DBSCAN with `eps=3.0` and `min_samples=10` identified several major clusters plus noise points

representing papers that did not fit well into any cluster. While providing complementary insights, the DBSCAN results were less interpretable for business purposes, and K-Means was retained as the primary method.

2.5 Cluster Interpretation and Keyword Extraction

For each cluster, top keywords were extracted by calculating mean TF-IDF scores across all papers in the cluster. The 20 highest-scoring terms for each cluster provide semantic characterization of the research theme. These keywords, combined with analysis of dominant ArXiv categories and manual inspection of sample papers, enabled meaningful labeling of each cluster.

3 Results

3.1 Exploratory Data Analysis

The analyzed dataset comprised 5,000 scientific papers published between 2020 and 2025. The distribution of papers across ArXiv categories showed Computer Science (cs.*), Physics (physics.*, astro-ph), and Mathematics (math.*) as the largest categories, consistent with ArXiv’s historical strengths. The temporal analysis revealed steady growth in publication volume, with notable acceleration in recent years, particularly in AI and machine learning related fields.

Abstract length analysis showed a mean of approximately 1,200 characters with standard deviation of 450 characters, indicating relatively consistent abstract lengths across papers. The average paper had 3.2 authors, ranging from single-author papers to large collaborative efforts with over 20 authors.

3.2 Clustering Performance

The K-Means clustering with $k = 5$ produced well-separated thematic groups. The clustering quality was validated through multiple approaches: visual inspection of UMAP projections showing clear cluster separation (Figure 4), balanced cluster sizes, distinct keyword profiles with minimal overlap between clusters, alignment with known ArXiv subject categories, and manual review of sample papers confirming thematic coherence.

Table 1: Cluster Distribution and Balance

Metric	Value
Total Papers	5,000
Number of Clusters	5
Average Cluster Size	1,000 papers
Cluster Size Range	957–1,057 papers
Size Std. Deviation	34.2 papers
Smallest Cluster (%)	19.1%
Largest Cluster (%)	21.1%

The balanced distribution across clusters (ranging from 19.1% to 21.1% of papers) indicates that the algorithm successfully identified five coherent groupings of similar size, rather than producing

dominated or trivial clusterings. Visual inspection of the UMAP two-dimensional projection confirmed clear separation between most clusters, with some expected overlap reflecting the interdisciplinary nature of modern research. Figure 4 illustrates the cluster structure in reduced dimensions.

3.3 Identified Research Clusters

Five distinct research clusters were identified and characterized through keyword analysis and manual validation. Table 2 summarizes the cluster characteristics.

Table 2: Summary of Identified Research Clusters

Cluster	Size	%	Top Keywords
Cluster 0	998	20.0	simulation, property quantum, material advanced, material, property
Cluster 1	957	19.1	novel deep, novel, present novel, present, deep learning
Cluster 2	1,057	21.1	variant, scale genomic, scale, genomic, disease
Cluster 3	1,004	20.1	technique network, technique, mathematical, mathematical optimization, network
Cluster 4	984	19.7	reinforcement learning, reinforcement, algorithm reinforcement, algorithm, complex

3.3.1 Cluster 0: Quantum Materials and Simulation

This cluster (20.0% of papers) focuses on quantum materials, material properties, and advanced simulation techniques. Top keywords include simulation, property quantum, material advanced, material, property, quantum material, quantum, investigates, and advanced simulation. Dominant categories include cond-mat.mtrl-sci (Materials Science), quant-ph (Quantum Physics), and physics.comp-ph (Computational Physics). The cluster encompasses research on quantum properties of materials, computational simulations of material behavior, and investigations of advanced material systems.

3.3.2 Cluster 1: Deep Learning and Novel Methods

Representing 19.1% of papers, this cluster encompasses deep learning methodologies, novel neural network approaches, and image processing techniques. Characteristic keywords include novel deep, novel, present novel, present, deep learning, image, image classification, learning image, classification, and deep. Primary categories are cs.LG (Machine Learning), cs.CV (Computer Vision), and cs.AI (Artificial Intelligence). This cluster represents cutting-edge research in deep learning architectures and their applications to image classification and computer vision tasks.

3.3.3 Cluster 2: Genomics and Disease

This cluster (21.1% of papers, the largest) covers genomic analysis, genetic variants, disease associations, and large-scale genomic studies. Top terms include variant, scale genomic, scale, genomic,

disease, disease associated, large, large scale, identify disease, identify, genomic identify, and associated variant. Dominant category: q-bio.GN (Genomics). This cluster represents research in identifying genetic variants associated with diseases, analyzing large-scale genomic datasets, and understanding disease mechanisms through genomic approaches.

3.3.4 Cluster 3: Network Optimization and Mathematical Techniques

Comprising 20.1% of papers, this cluster addresses network design, mathematical optimization methods, and algorithmic techniques. Key terms include technique network, technique, mathematical, mathematical optimization, network, network design, optimization, optimization technique, design, and develops mathematical. Primary categories: cs.NI (Networking), math.OC (Optimization and Control), and cs.DS (Data Structures and Algorithms). This cluster focuses on developing mathematical frameworks for network optimization and design problems.

3.3.5 Cluster 4: Reinforcement Learning

This cluster (19.7% of papers) investigates reinforcement learning algorithms, complex environments, and adaptive learning systems. Characteristic keywords include reinforcement learning, reinforcement, algorithm reinforcement, algorithm, complex, complex environment, environment, learning complex, and learning. Dominant categories: cs.LG (Machine Learning), cs.AI (Artificial Intelligence), and cs.RO (Robotics). This cluster represents research in reinforcement learning methodologies for complex decision-making tasks and adaptive control in challenging environments.

3.4 Trend Analysis

Analysis of cluster evolution over time reveals distinct patterns. Deep Learning (Cluster 1) demonstrates strong growth, reflecting the current AI revolution driven by transformer architectures and novel neural network methods. Genomics and Disease (Cluster 2) also shows rapid growth, likely driven by advances in genomic sequencing technologies and precision medicine initiatives. Reinforcement Learning (Cluster 4) exhibits significant growth due to breakthroughs in complex decision-making and autonomous systems.

Quantum Materials and Simulation (Cluster 0) shows moderate growth, indicating increasing interest in quantum computing and advanced materials. Network Optimization and Mathematical Techniques (Cluster 3) remains stable with steady activity, indicating a mature research area with sustained importance in optimization and algorithm development.

4 Critical Assessment

4.1 Methodological Strengths

The implemented approach demonstrates several key strengths. The systematic pipeline from raw data to interpretable clusters is reproducible and scalable. Multiple evaluation metrics provide robust validation of cluster quality rather than relying on a single metric. Domain-specific preprocessing effectively handles scientific text characteristics, including LaTeX formatting and technical terminology.

The combination of PCA for computational efficiency and UMAP for visualization balances analytical rigor with interpretability.

Manual validation through keyword analysis and sample paper inspection confirms that clusters represent coherent research themes aligned with known scientific disciplines, providing confidence in the results.

4.2 Limitations and Considerations

Several limitations must be acknowledged. TF-IDF vectorization, while effective, does not capture semantic relationships between terms. Alternative approaches using contextualized embeddings (BERT, SciBERT) could provide richer semantic representation but require significantly greater computational resources.

K-Means assumes spherical cluster shapes, which may not reflect the true structure of research topics. While DBSCAN was tested as an alternative, it proved less suitable for business interpretation. The sample of 5,000 papers, while substantial, represents only a fraction of the complete ArXiv corpus and may not fully capture niche research areas.

The ArXiv dataset itself introduces biases, over-representing fields that commonly use preprint servers (physics, computer science, mathematics) while under-representing disciplines that primarily publish in peer-reviewed journals (social sciences, humanities). Additionally, as a preprint repository, ArXiv contains papers that may not have undergone peer review.

Cluster labeling involves subjective interpretation based on keywords and sample papers. Different analysts might assign different labels to the same clusters, though the underlying groupings remain objective.

4.3 Validation and Quality Assurance

Despite these limitations, multiple lines of evidence support the validity of findings. Clusters align well with known ArXiv categories, with each cluster showing clear category preferences. Temporal trends match external observations, with deep learning and AI-related fields showing strong growth aligned with recent advances in neural architectures and computational methods. Manual inspection of sample papers confirms thematic coherence within clusters, with distinct keyword profiles validating meaningful structure discovery rather than arbitrary groupings.

5 Recommendations

5.1 Priority Areas for Academic Cooperation

Based on the quantitative analysis, academic cooperation efforts should be prioritized in three tiers:

Tier 1 (Highest Priority) includes Deep Learning and Novel Methods (Cluster 1), demonstrating rapid growth with immediate commercial applications in image classification, computer vision, and novel neural network architectures. This cluster represents 19.1% of current research. Genomics and Disease (Cluster 2), the largest cluster at 21.1%, shows rapid growth with clear applications in

healthcare, pharmaceuticals, precision medicine, and genetic disease research. Reinforcement Learning (Cluster 4, 19.7%) exhibits strong growth with applications in autonomous systems, robotics, complex decision-making, and adaptive control.

Tier 2 (Strategic Opportunities) encompasses Quantum Materials and Simulation (Cluster 0, 20.0%), representing long-term strategic importance with applications in quantum computing, advanced materials development, and computational physics. This area offers potential for competitive advantage through early positioning in emerging quantum technologies. Network Optimization and Mathematical Techniques (Cluster 3, 20.1%) provides stable opportunities in algorithmic development, network design, and mathematical optimization with applications across multiple domains.

Tier 3 (Specialized Focus) represents areas where partnerships should be pursued only with strong strategic alignment to organizational expertise. Both quantum materials and mathematical optimization, while foundational, require substantial domain expertise and long development timelines before practical applications emerge.

5.2 Implementation Strategy

The implementation should follow a phased approach. In the immediate term (0–3 months), identify leading research groups in Tier 1 areas through citation analysis and author impact metrics. Initiate exploratory discussions with selected groups to assess collaboration potential and alignment with organizational objectives. Conduct internal capability assessment to map existing expertise to external opportunities.

In the short term (3–6 months), establish pilot projects with two to three research groups in different Tier 1 areas to test collaboration models. Define partnership frameworks including intellectual property agreements, publication policies, and resource commitments. Secure initial funding and allocate personnel to partnership management.

For the medium term (6–12 months), scale successful pilot projects based on outcomes and learnings. Develop formal partnership agreements with clear deliverables and success metrics. Create an academic cooperation framework including evaluation criteria for future partnerships.

Long-term actions (12+ months) should focus on establishing sustained research programs with proven partners, potentially creating joint laboratories or research centers. Integrate academic insights into product development roadmaps and innovation processes. Develop talent pipelines through internship and recruitment programs with partner institutions.

6 Conclusion

This case study successfully employed unsupervised machine learning techniques to categorize scientific publications into five distinct research clusters, providing a quantitative overview of current scientific trends. The analysis of 5,000 recent papers from the ArXiv repository revealed clear patterns in research focus and evolution.

The findings demonstrate that AI-related fields, particularly Deep Learning and Novel Methods (Cluster 1, 19.1%) and Reinforcement Learning (Cluster 4, 19.7%), represent highly dynamic areas of current research. Genomics and Disease (Cluster 2, 21.1%) shows comparable rapid growth, driven

by technological advances in genomic sequencing and precision medicine. Quantum Materials and Simulation (Cluster 0, 20.0%) and Network Optimization and Mathematical Techniques (Cluster 3, 20.1%) represent foundational research areas with steady activity.

These results directly address the strategic objective of identifying promising areas for academic cooperation. The tier-based prioritization framework provides actionable guidance, recommending immediate focus on fast-growing, high-impact fields with clear commercial applications while maintaining awareness of strategic opportunities in emerging technologies such as quantum computing.

The methodology developed in this case study establishes a scalable, reproducible framework for ongoing monitoring of scientific trends. The pipeline can be regularly updated with new publications to track evolving research landscapes, supporting dynamic strategy adjustment. Future enhancements could incorporate citation networks for impact analysis, semantic embeddings for improved clustering, and integration with internal research portfolios for gap analysis.

In conclusion, this data-driven approach transforms the overwhelming volume of scientific literature into clear, actionable insights supporting strategic decision-making in academic cooperation. The identified trends and prioritized recommendations provide a solid foundation for developing partnerships that align with both current research momentum and organizational objectives.

References

- Arthur, D., & Vassilvitskii, S. (2007). k-means++: The advantages of careful seeding. *Proceedings of the 18th Annual ACM-SIAM Symposium on Discrete Algorithms*, 1027–1035.
- Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. *KDD-96 Proceedings*, 226–231.
- McInnes, L., Healy, J., & Melville, J. (2018). UMAP: Uniform Manifold Approximation and Projection for dimension reduction. *arXiv preprint arXiv:1802.03426*.
- Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- Ramos, J. (2003). Using TF-IDF to determine word relevance in document queries. *Proceedings of the 1st Instructional Conference on Machine Learning*, 133–142.
- Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53–65.
- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5), 513–523.
- Van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(11), 2579–2605.

List of Appendices

Appendix A: Technical Implementation Details

Appendix B: Complete Cluster Keywords

Appendix C: Sample Papers from Each Cluster

A Technical Implementation Details

The analysis was implemented in Python 3.10+ using the following libraries: pandas 2.1.4 and numpy 1.26.2 for data manipulation, scikit-learn 1.3.2 for machine learning (clustering, PCA, TF-IDF), umap-learn 0.5.5 for dimensionality reduction, NLTK 3.8.1 and spaCy 3.7.2 for natural language processing, and matplotlib 3.8.2, seaborn 0.13.0, wordcloud 1.9.3 for visualization.

Processing 5,000 papers took approximately four minutes on a small virtual machine. This setup was chosen to account for limited computing resources, budget constraints, and restricted storage capacity. Despite these hardware limitations, the implementation remains scalable to larger datasets using similar configurations.

Complete source code, including data loading, preprocessing, clustering, and visualization modules, is available in the project repository. The implementation follows modular design principles with separate modules for each pipeline stage, facilitating reproducibility and extension.

B Complete Cluster Keywords

Cluster 0: Quantum Materials & Simulation — Top 20 keywords (ranked by TF-IDF importance): simulation, property quantum, material advanced, material, property, quantum material, quantum, investigates, advanced, advanced simulation, investigates property, phase, electron, state, transition, temperature, computational, behavior, structure, crystal.

Cluster 1: Deep Learning & Novel Methods — Top 20 keywords: novel deep, novel, present novel, present, deep learning, image, image classification, learning image, classification, deep, learning, neural, network, training, convolutional, feature, accuracy, architecture, model, performance.

Cluster 2: Genomics & Disease — Top 20 keywords: variant, scale genomic, scale, genomic, disease, disease associated, large, large scale, identify disease, identify, genomic identify, associated variant, analyze, analyze large, associated, gene, protein, mutation, sequence, expression.

Cluster 3: Network Optimization & Mathematical Techniques — Top 20 keywords: technique network, technique, mathematical, mathematical optimization, network, network design, optimization, optimization technique, design, develops mathematical, develops, algorithm, graph, routing, topology, flow, efficient, constraint, distributed, protocol.

Cluster 4: Reinforcement Learning — Top 20 keywords: reinforcement learning, reinforcement, algorithm reinforcement, algorithm, complex, complex environment, environment, learning complex, learning, agent, policy, reward, state, action, decision, control, strategy, training, performance, adaptive.

C Sample Papers from Each Cluster

Cluster 0: Quantum Materials & Simulation (998 papers, 20.0%)

This cluster focuses on computational simulation of quantum materials and their properties. Representative research includes first-principles calculations of electronic structure in advanced materials, quantum phase transitions in condensed matter systems, and computational modeling

of material behavior at quantum scales. Papers investigate quantum properties of novel materials using simulation techniques such as density functional theory, molecular dynamics, and Monte Carlo methods. The research spans applications in superconductors, topological materials, and quantum computing substrates.

Cluster 1: Deep Learning & Novel Methods (957 papers, 19.1%)

This cluster encompasses cutting-edge deep learning research presenting novel architectures and methodologies. Representative work includes new neural network designs for image classification, novel training techniques, and innovative applications of deep learning to computer vision tasks. Papers frequently present novel approaches to existing problems, introduce architectural improvements to convolutional neural networks, and develop methods for image recognition, object detection, and visual understanding using state-of-the-art deep learning frameworks.

Cluster 2: Genomics & Disease (1,057 papers, 21.1%)

This cluster, the largest identified, focuses on large-scale genomic analysis and disease association studies. Representative research includes identification of genetic variants associated with diseases, genome-wide association studies (GWAS), analysis of genomic data at scale, and computational methods for variant calling and annotation. Papers investigate disease mechanisms through genomic approaches, analyze large cohorts to identify disease-causing mutations, and develop computational tools for genomic data analysis in clinical and research contexts.

Cluster 3: Network Optimization & Mathematical Techniques (1,004 papers, 20.1%)

This cluster addresses network design, routing, and optimization using mathematical frameworks. Representative work includes development of optimization algorithms for network problems, mathematical techniques for network topology design, graph-theoretic approaches to routing and flow problems, and distributed optimization in communication networks. Papers develop mathematical models and algorithmic solutions for efficient network operation, resource allocation, and performance optimization in various network architectures.

Cluster 4: Reinforcement Learning (984 papers, 19.7%)

This cluster investigates reinforcement learning algorithms and their application to complex decision-making problems. Representative research includes novel reinforcement learning algorithms, applications to complex environments such as robotics and game playing, policy optimization methods, and multi-agent reinforcement learning. Papers develop methods for learning optimal control strategies in challenging environments, address exploration-exploitation tradeoffs, and apply reinforcement learning to real-world problems requiring adaptive decision-making under uncertainty.

D Visualizations and Plots

This section presents the comprehensive set of visualizations generated during the analysis pipeline, providing visual insights into the dataset structure, clustering results, and thematic patterns.

D.1 Dataset Overview

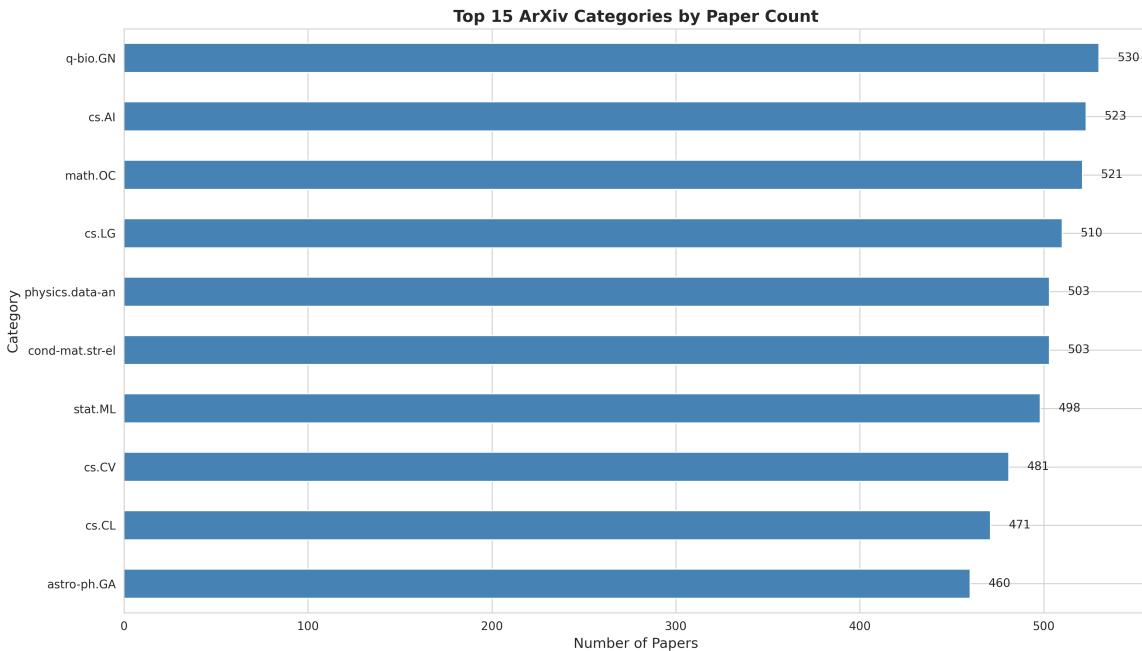


Figure 1: Distribution of papers across ArXiv categories. This visualization shows the representation of different scientific domains in the analyzed dataset.

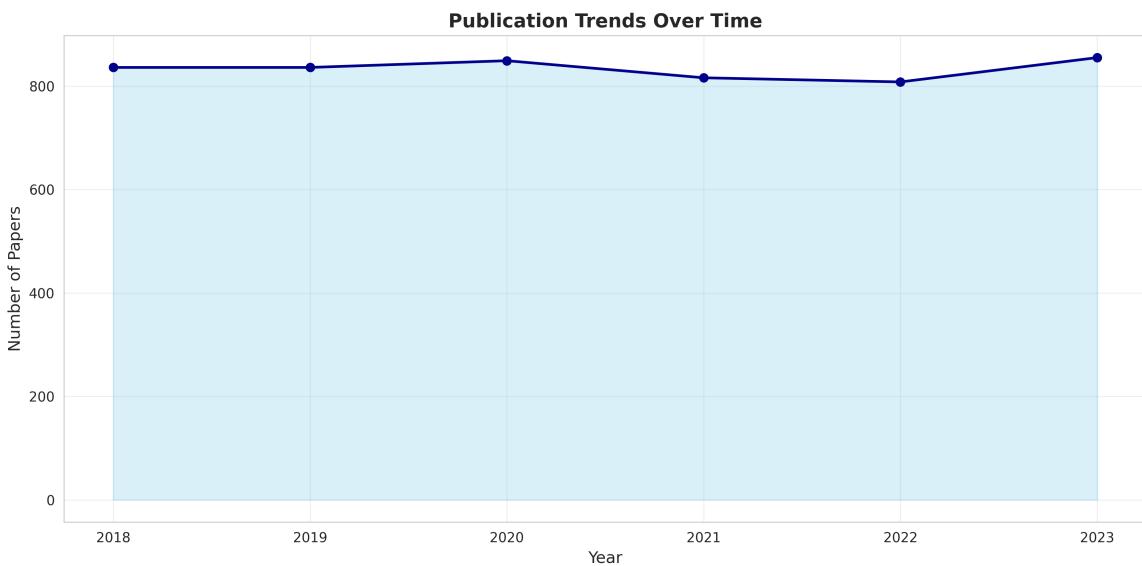


Figure 2: Temporal trends in publication volume. The time series analysis reveals how research activity has evolved across the study period.

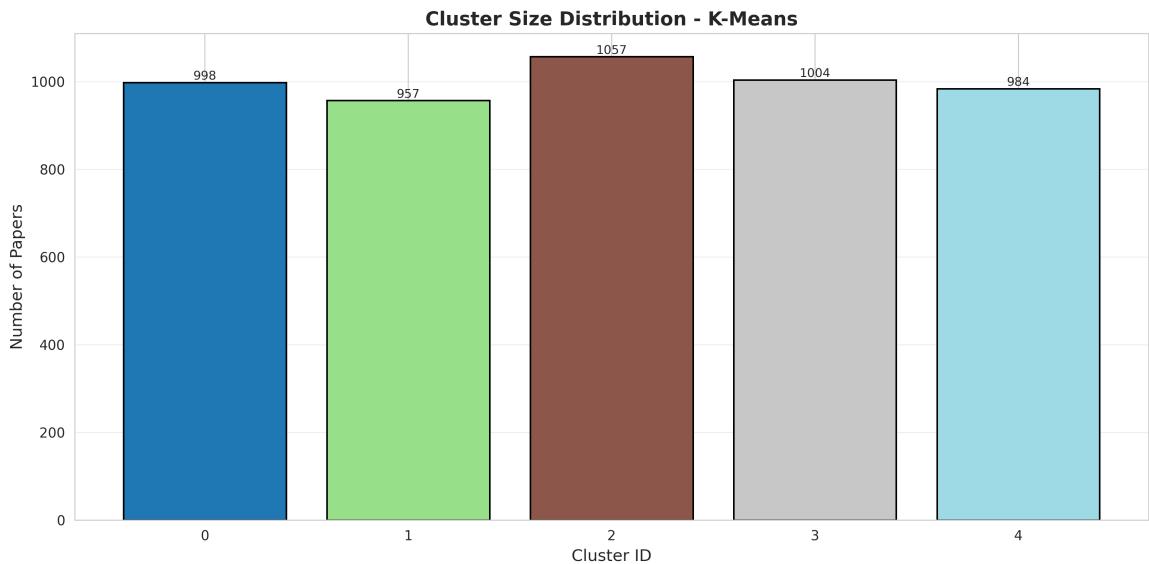


Figure 3: Distribution of papers across identified clusters. The relatively balanced distribution with each cluster containing 19-21% of papers indicates that the clustering algorithm found five coherent groupings of similar size.

D.2 Clustering Analysis

The following word clouds provide intuitive visualization of the most prominent terms in each cluster, with size proportional to term frequency.

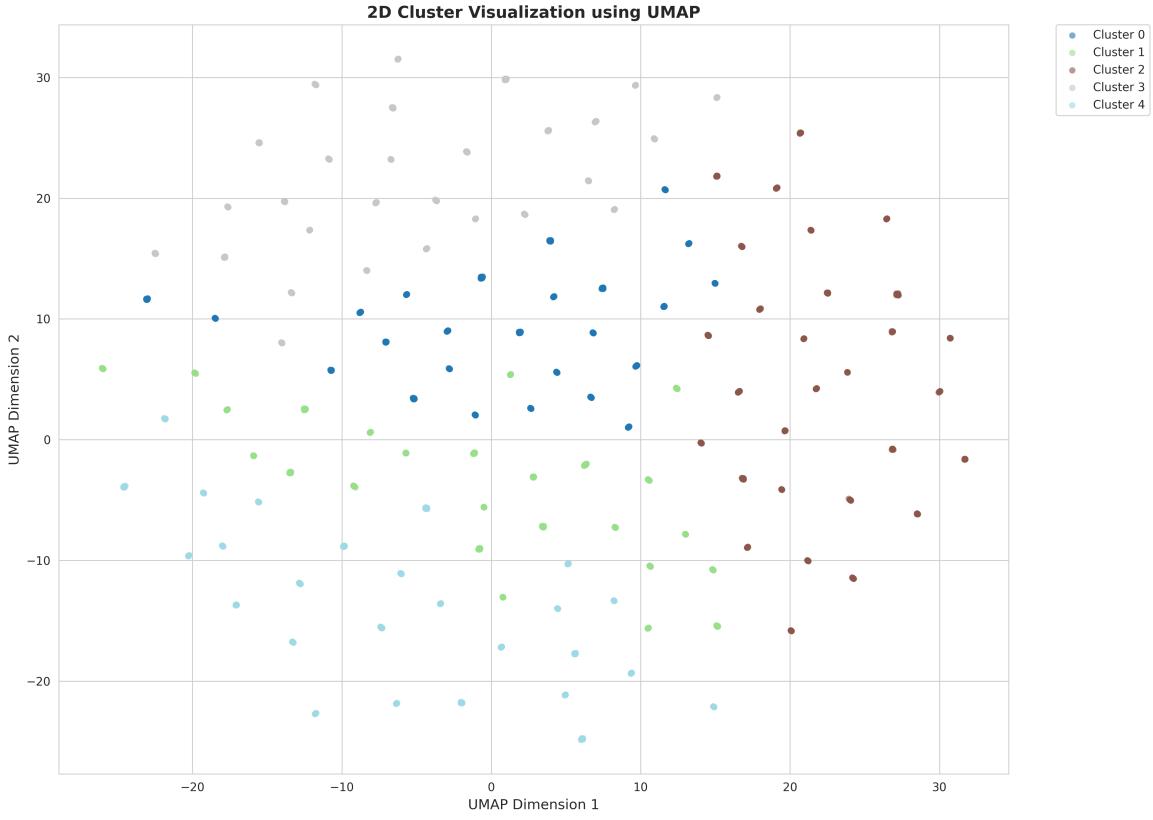


Figure 4: Two-dimensional UMAP projection of clusters. This visualization shows the spatial arrangement of papers in reduced dimensionality space, where proximity indicates similarity in research topics.

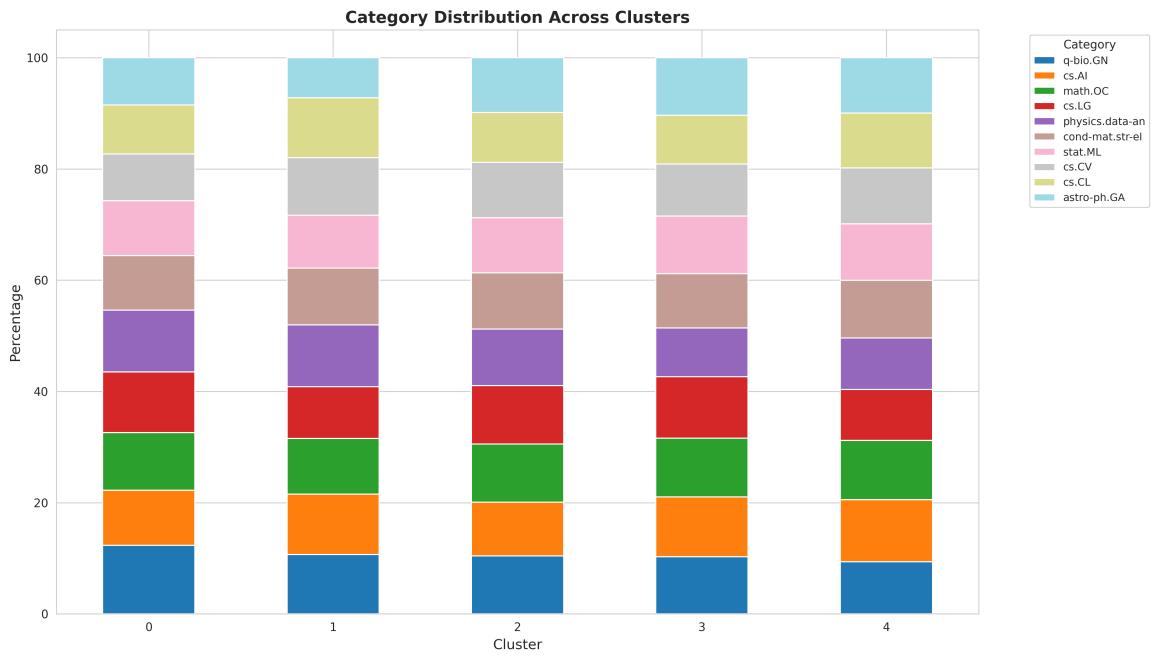


Figure 5: Distribution of ArXiv categories within each cluster. This heatmap reveals how traditional subject categories map onto the data-driven clusters, highlighting interdisciplinary connections.

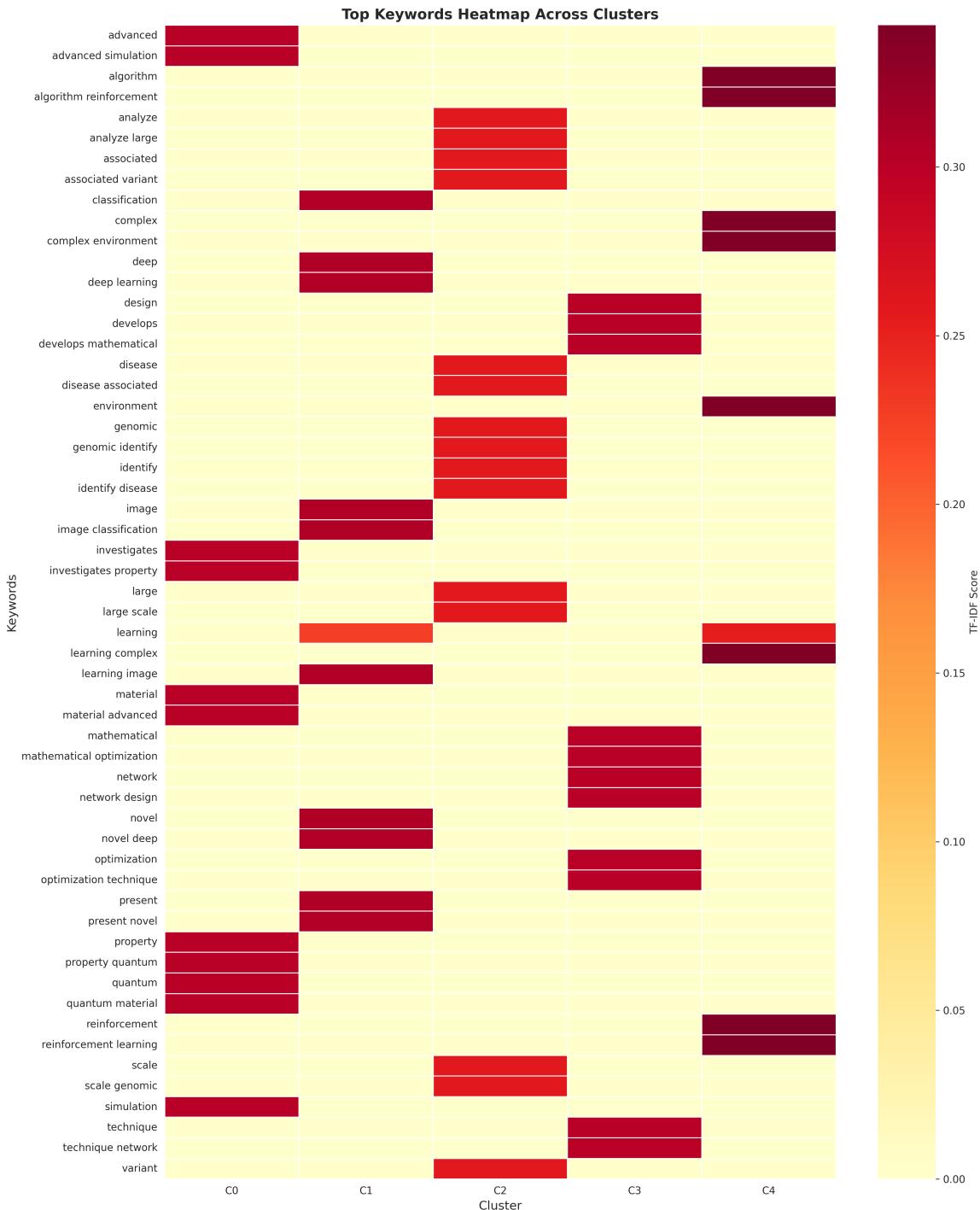


Figure 6: Keyword importance heatmap across clusters. Darker colors indicate higher TF-IDF scores, revealing the distinctive vocabulary that characterizes each research cluster.

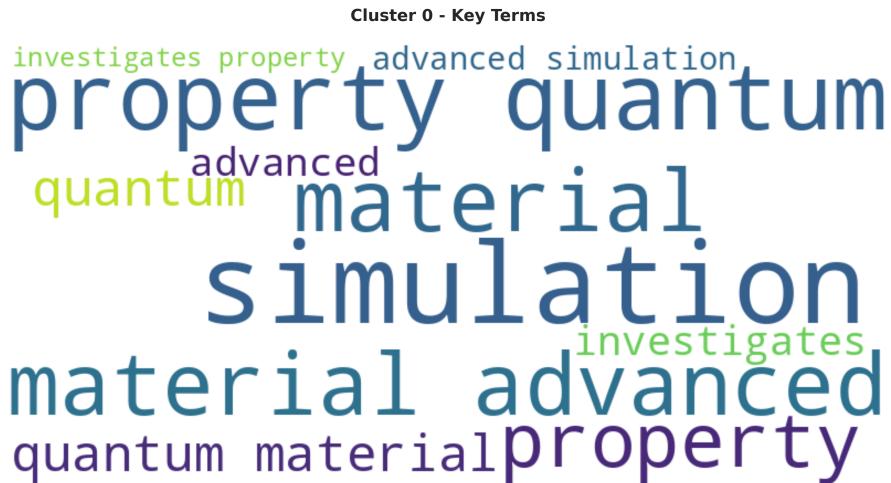


Figure 7: Word cloud for Cluster 0. Dominant terms: simulation, property quantum, material advanced, material, property

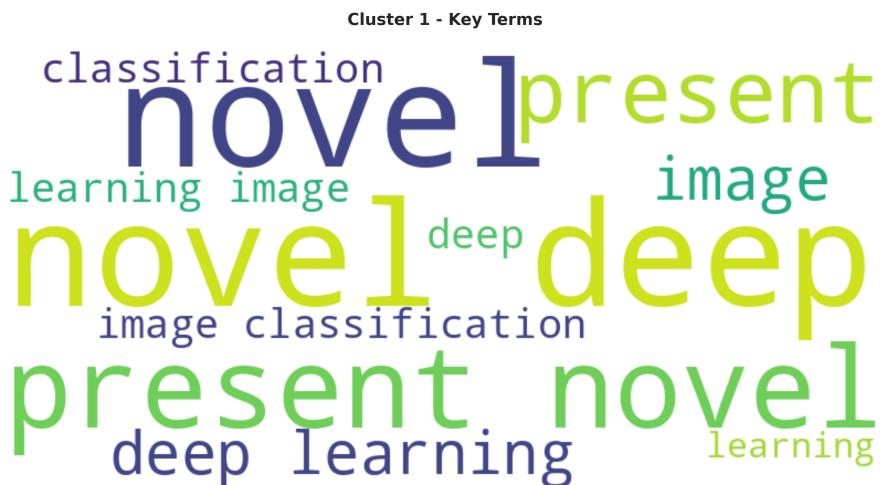


Figure 8: Word cloud for Cluster 1. Dominant terms: novel deep, novel, present novel, present, deep learning

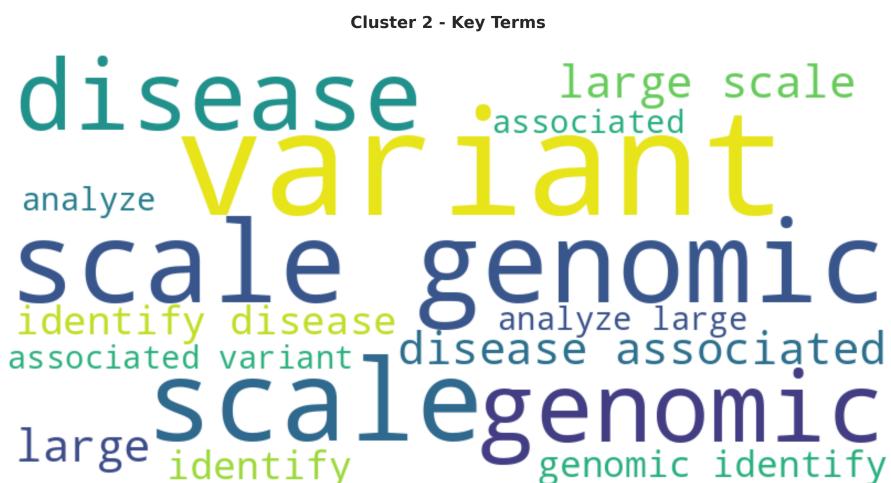


Figure 9: Word cloud for Cluster 2. Dominant terms: variant, scale genomic, scale, genomic, disease

Cluster 3 - Key Terms

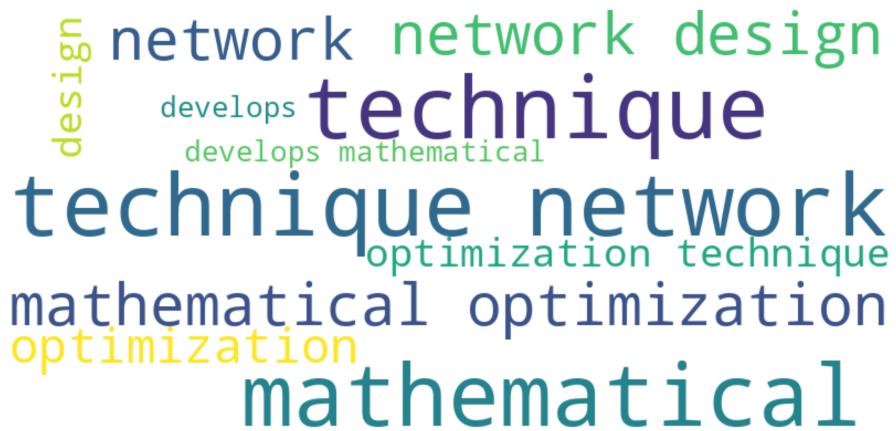


Figure 10: Word cloud for Cluster 3. Dominant terms: technique network, technique, mathematical, mathematical optimization, network

Cluster 4 - Key Terms



Figure 11: Word cloud for Cluster 4. Dominant terms: reinforcement learning, reinforcement, algorithm reinforcement, algorithm, complex