

Intelligent Recommendation and Indexing System for Deep Learning Papers

1st Sihan Chen

Fu Foundation School of Engineering and Applied Science
Columbia University
New York, United States
sc5534@columbia.edu

2nd Haowei Chen

Fu Foundation School of Engineering and Applied Science
Columbia University
New York, United States
hc3515@columbia.edu

Abstract—This paper introduces a personalized learning platform that integrates a domain-specific knowledge graph, semantic embeddings, and retrieval-augmented generation (RAG) to streamline the exploration of deep learning literature. The platform constructed a heterogeneous knowledge graph to represent relationships between research papers, methods, and tasks. We use Graph Attention Network (GAT) to predict user-paper interactions, and we achieved an AUC of 0.91 in recommendation tasks. Our system supports vectorized search, paper summarization, and open-ended question answering. Our platform has well-designed front-end surface to show the platform’s effectiveness in simplifying literature discovery and providing personalized insights.

I. INTRODUCTION

The advancements in artificial intelligence (AI) and deep learning (DL) have changed various domains, like natural language processing, computer vision, biomedical research, etc. The emergence of technologies like AlphaGo that defeated the world Go champion [1]. ChatGPT has transformed daily life [2] and illustrates the rapid progress of AI. Reflecting this growth, the number of published DL-related papers rose dramatically from 5,000 in 2016 to over 15,000 by 2019 [3]. This widespread interest is fueled by the versatility of neural networks. It enables applications in diverse areas such as computational biology, network security, and video generation.

While the interdisciplinary nature of DL draws broad attention, it also presents challenges. The rapid proliferation of research papers makes readers difficult to understand key concepts. DL algorithms require strong mathematical foundations, but their intricate neuron-level processes often confuse beginners. Furthermore, the field’s cross-disciplinary nature and dense citation networks compound these challenges. Essential innovations are frequently obscured in paper titles; for example, the Retinanet paper emphasizes “Focal Loss” [4], while the GPT-3 paper highlights “Few-shot Learning” [2]. Conflicts between works, such as critiques of Mamba by Mamba Out [5], [6], further complicate literature navigation.

Addressing these issues requires a system that can simplify the process of finding and understanding research. We believe such a system should be objective, user-focused, and tailored to individual preferences. This allows both learners of different level to explore the DL research landscape efficiently.

This project introduces a personalized learning platform combining a DL knowledge graph, retrieval-augmented generation (RAG), and graph neural networks (GNNs). Our system aims to bridge the gap between the overwhelming volume of information and users’ unique needs by offering tailored paper recommendations and efficient retrieval tools.

The platform contributes the following innovations:

- A scalable DL knowledge graph capturing semantic relationships such as citations, methods, and application areas.
- Enhanced paper search and summarization through LLM embeddings and RAG. They integrate structured graph data and external sources.
- A recommendation system predicting user-paper interactions to provide personalized suggestions.

This report details the system’s background, design, and implementation. We’ll show its ability to streamline academic exploration and improve the learning experience for various users. Specifically, the performance of the two-layer Graph Attention Network (GAT) was rigorously evaluated. The results demonstrate the model’s ability to accurately capture citation relationships and semantic connections within the knowledge graph, with an AUC of 0.91 in prediction and recommendation. Our system offers precise recommendations and generating detailed, context-aware answers to user queries.

II. RELATED WORK

A. Knowledge Graphs

Knowledge graphs (KGs) are structured knowledge bases that represent information through interconnected entities and relationships. Formally, a KG can be expressed as $G = \{E, R, T\}$, where E denotes the set of entities, R denotes the set of relationships, and T represents the ontology defining connections between entities and relationships:

$$T = \{t_1, t_2, \dots, t_N\}, \quad t_i = (e_j, r_i, e_k), \quad e_j, e_k \in E, \quad r_i \in R.$$

The concept of knowledge graphs was first introduced in 1972 by Edgar W. Schneider in discussions on modular course guidance systems [7]. Later developments include DBpedia and Freebase in 2007, which leveraged user-generated content to build general-purpose KGs. Google redefined the concept

in 2012 and made it more versatile [8]. Wolfram researchers formalized its mathematical foundations in 2016 [9]. Knowledge graphs are now widely applied in search engines and question-answering systems.

In academic domains, citation networks act as specialized KGs, with papers and their citation relationships forming the primary entities and links. Examples include the Cora dataset (2,708 papers and 5,429 citation links) [10], Citeseer dataset (3,312 papers and 4,732 citation links) [11], and ArXiv HEP-TH dataset (27,770 papers and 352,807 citation links) [12]. These datasets are critical for tasks like topic modeling, text classification, and summarization. Due to their dense connectivity and domain-specific focus, citation networks in KGs provide high knowledge density and strong guidance for exploring specific research areas.

B. Graph Neural Networks (GNNs)

Graph neural networks (GNNs) are specialized neural networks that directly operate on graph structures. GNNs iteratively update parameters and transmit information:

$$H_{t+1} = f(H_t, x),$$

where H represents the graph's node embeddings, x denotes node features, and f is the updating function.

Classic GNN models include Graph Convolutional Networks (GCN) and Graph Attention Networks (GAT), etc. They balance computational complexity and accuracy by limiting the scope of neighboring nodes during updates. GNNs are widely used in edge classification, node classification, and graph representation learning. They have wide application scenarios spanning biology, social networks, recommendation systems, and cybersecurity [13].

GNNs also enhance KG-based recommendation systems. For instance, KGAT employs attention mechanisms to train vector representations for users and items. It achieved state-of-the-art performance with low computational costs [14]. Similarly, CGAT combines local and global information to extract critical subgraph features, and demonstrated the power of GNNs in integrating user interactions and KG relationships for effective recommendations [15].

C. Large Language Models (LLMs)

Large language models (LLMs) developed quickly from the transformer architecture introduced in "Attention Is All You Need" [16]. LLMs represent a paradigm shift in AI. Since 2018, OpenAI's GPT series has set benchmarks for generative tasks. Later developments like LLaMA and Mixtral emerged as open-source alternatives [17], [18]. LLMs typically feature billions of parameters, which makes full retraining impractical. Instead, their usage focuses on fine-tuning and prompt engineering.

LLMs are known for their few-shot learning capabilities. This enables tasks with minimal training examples [19]. However, a critical challenge is hallucination, where models produce inaccurate or irrelevant outputs. These errors are usually

categorized into factual hallucinations (involving incorrect information) and faithfulness hallucinations (outputs inconsistent with input queries) [20]. Contributing factors include biased datasets, suboptimal usage of data, and insufficient attention mechanisms.

D. Retrieval-Augmented Generation (RAG)

Retrieval-augmented generation (RAG) can mitigate LLM hallucinations by combining model-generated responses with external knowledge bases. RAG retrieves relevant external information during the generation process. They can provide accurate and explainable answers in open-domain question answering tasks [21]. The integration of retrieval and generation forms a robust framework for handling knowledge-intensive tasks. As a result, contextual relevance and factual accuracy are increased.

RAG typically involves two main components: a retriever and a generator. Given a query, the retriever searches an external knowledge corpus (e.g. a knowledge graph) to retrieve relevant documents or facts. The generator, in our case an LLM, synthesizes a response by combining retrieved information with its own generative capabilities. Our structure can improve the factual reliability of structured or unstructured external data and the natural language fluency of the LLM.

Despite significant advances in these fields, existing solutions often fail to seamlessly integrate these technologies for personalized learning platforms. For instance, traditional knowledge graphs like DBpedia and Freebase are good at organizing general information. However, their application-specific adaptation is limited. On the other hand, state-of-the-art LLMs can provide generative capabilities. But they are struggling with factual consistency and domain adaptation without external augmentation.

Our approach bridges these gaps by integrating an RAG pipeline with a domain-specific knowledge graph enhanced by GNNs. Compared to conventional systems, our platform dynamically combines structured graph queries and personalized recommendations. Our goal is to address the unique challenges of navigating dense and cross-disciplinary research fields.

III. DATASET

Our dataset consists of a diverse range of structured and unstructured data relevant to deep learning research. The primary data sources include ArXiv, Papers With Code, and Semantic Scholar. They contribute to the complementary aspects of the dataset. Figure 1 illustrates each step from data extraction and preprocessing to storage in structured databases. ArXiv provides bibliographic information such as titles, abstracts, and publication dates. Papers With Code focuses on task and method-specific metadata. Semantic Scholar offers citation relationships and additional bibliometric data.

The dataset creation process involves multiple stages. It begins with data acquisition. ArXiv metadata is retrieved via its publicly available API. They usually have unique identifiers (e.g., 1512.03385) to track publication details. Similarly,

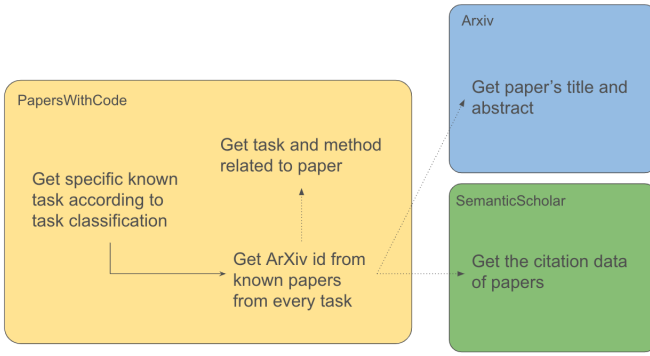


Fig. 1. Our data acquisition pipeline.

Papers With Code provides a structured taxonomy of research domains and subdomains. This helps categorizing research papers by tasks and methods. Semantic Scholar is queried to extract citation relationships and enrich metadata, and ensures comprehensive bibliometric coverage.

Preprocessing ensures data quality and consistency. Titles and abstracts are tokenized and cleaned to remove anomalies and standardize formatting. Metadata from multiple sources is merged. We resolved duplicate entries through identifier matching and manual verification where necessary. Semantic Scholar’s API is employed to validate citation relationships, which increases the accuracy and completeness of the graph representation.

The processed data is transformed into a heterogeneous knowledge graph stored in a Neo4j database. Nodes in the graph represent papers, tasks, and methods, while edges capture relationships such as citations, applications, and usage. This graph-based structure enables complex queries and facilitates downstream tasks such as personalized recommendations and retrieval-augmented generation.

We employ asynchronous processing techniques to optimize efficiency. A coroutine-based mechanism handles API requests and balances retrieval speed with server load. The curated dataset comprises over 30,000 papers, 10,000 citation relationships, and metadata spanning more than 100 research tasks and methods.

IV. METHOD

This section outlines the methods used to address the challenges introduced in the study. These approaches were selected based on their ability to effectively utilize the rich, heterogeneous data required for deep learning research exploration.

A. Co-routine Data Preprocessing

To handle the large-scale asynchronous retrieval of data, a co-routine pool was implemented. This mechanism optimizes API requests, ensuring parallel execution without exceeding server rate limits. The collected data undergoes a rigorous preprocessing phase to clean and standardize information. Titles and abstracts are tokenized, inconsistencies are removed, and

citation relationships are verified against Semantic Scholar’s API.

B. Knowledge Graph Construction

A heterogeneous knowledge graph is constructed to represent the relationships. Figure 2 illustrates the ER diagram for the system. Nodes in the knowledge graph correspond to entities such as research papers, machine learning methods, and application tasks. Edges capture relationships like citations, methodological applications, and task associations. The graph structure is designed to support both structural queries and algorithmic learning.

To balance expressiveness and scalability, the knowledge graph is stored in Neo4j, while auxiliary data (e.g., user preferences) is managed in a MySQL database. Our dual-database approach enables complex graph queries alongside efficient handling of structured data. The graph’s modular design supports integration with algorithms such as GNNs for recommendation and RAG for enhanced query responses.

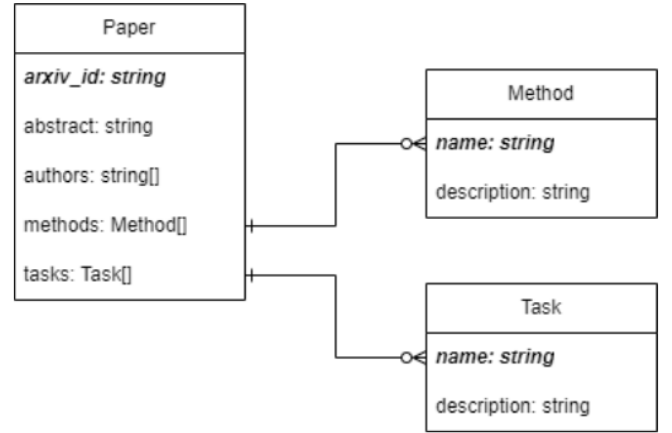


Fig. 2. The Entity-Relationship Diagram of our data pipeline.

C. Vector Retrieval

The vector retrieval functionality ensures efficient and accurate searches for relevant documents based on semantic similarity. The system transforms text into high-dimensional vectors.

The retrieval process begins with embedding generation. Each document, such as a research paper’s title, abstract, or full text, is encoded into a vector representation. Pre-trained language model Google Gemini is fine-tuned to produce embeddings that reflect the contextual relevance of the text. These embeddings are then stored in a dedicated vector database. They form the basis of the retrieval engine. To maintain accuracy and adaptability, the platform periodically updates embeddings in order to make sure that newly added documents are indexed and retrievable. We aim to keep the system aligned with the latest research trends.

To find relevant documents, the system performs similarity searches by comparing query embeddings to those in the

database. Similarity is typically measured using cosine similarity or Euclidean distance. The platform employs optimized vector databases of Weaviate, which utilize approximate nearest neighbor (ANN) search algorithms for high efficiency. ANN techniques significantly reduce computational overhead. It supports real-time query responses even with large datasets. Moreover, the system supports multi-modal queries, where users can search using text inputs (e.g., keywords or abstracts) or structural elements (e.g., node embeddings from the knowledge graph).

A defining feature of this functionality is its integration with the knowledge graph. Once relevant documents are retrieved, their corresponding nodes in the graph are identified. This process enhances the retrieval process by providing additional relational insights. For example, a query on "Graph Neural Networks for Recommendation Systems" might retrieve both foundational and application-focused papers. The system further enriches results by highlighting citation chains, co-authorship networks, or shared methodologies.

To ensure scalability, the vector database is partitioned into shards based on thematic domains or publication years. This can reduce the search space and improving retrieval speeds. Frequently accessed queries are handled efficiently through caching mechanisms.

D. Retrieval-Augmented Generation (RAG)

RAG combines document retrieval with language model generation to address the limitations of large language models in knowledge-intensive tasks. The retrieval component identifies relevant documents using vector embeddings generated by Google Gemini. The selected documents are then passed to a language model, which synthesizes a response that integrates retrieved knowledge with generative capabilities.

RAG excels in answering open-ended queries by grounding generated outputs in factual data. For instance, a user query like "Explain U-Net" triggers the retrieval of relevant papers and produces a comprehensive response describing the model's architecture and applications. Compared to standalone language models, RAG improves accuracy and context-awareness while mitigating issues like hallucinations.

E. Recommendation System

The recommendation system is powered by a Graph Attention Network (GAT), designed to predict user-paper interactions and provide personalized suggestions. GAT employs attention mechanisms to weigh the importance of neighboring nodes during embedding updates. The core update rule is defined as:

$$H_i^{(l+1)} = \text{ReLU} \left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} W^{(l)} H_j^{(l)} \right),$$

where α_{ij} represents the attention coefficients calculated based on the features of connected nodes. These coefficients help the network selectively focus on important connections, enhancing its ability to capture meaningful relationships.

Building upon this, the GAT-based system uses data from the deep learning knowledge graph, enriched with semantic embeddings generated from large language models. The integration of multi-hop citation relationships and text-based embeddings allows the recommendation system to enhance both the relevance and diversity of its suggestions.

The nodes represent papers and users, and edges denote interactions or citation relationships. Semantic embeddings derived from external sources are incorporated to enhance node features. The multi-layer GAT dynamically updates node representations through iterative attention-weighted aggregations. This ensures the final embeddings are both context-aware and representative of the user's preferences. Finally, A binary classification task is used to predict the likelihood of user-paper interactions. The system ranks papers based on these predictions to generate a personalized recommendation list.

The GAT-based approach was chosen over traditional collaborative filtering methods because it effectively captures the graph-structured nature of the data and integrates node-level features. Initial attempts to employ simpler methods such as matrix factorization and item-based collaborative filtering resulted in less precise recommendations. They often underscore the importance of utilizing graph-based insights.

F. Alternative Approaches

Alternative methods were considered during the development phase. Simpler approaches, such as keyword-based search or traditional recommendation systems, were tested but unable to achieve the desired level of personalization and contextual relevance. For instance, keyword matching often returned unrelated results due to the ambiguity of terms like "GPT." Similarly, collaborative filtering models struggled to adapt to the dynamic and interdisciplinary nature of deep learning research.

The selected approaches address these shortcomings by combining graph-based learning, semantic embeddings, and generative models. This integration provides a comprehensive framework for navigating complex research landscapes.

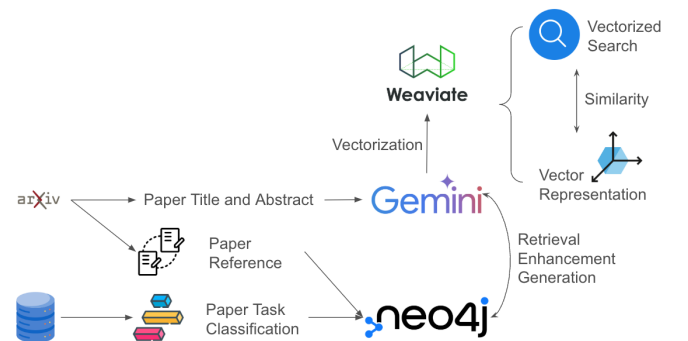


Fig. 3. Data Transmission Map: A visual representation of data flow within the platform.

V. SYSTEM OVERVIEW

A. System Architecture

The personalized learning platform is designed with a modular and scalable architecture to ensure efficient interaction and robust recommendations for diverse user needs. Figure 3 illustrates the data transmission architecture. It details the flow and integration of various components.

Key components of the system include:

1. Knowledge Graph Storage: The Neo4j database manages the graph-structured data. It supports operations like entity management and multi-hop citation queries through the Cypher query language. For structured queries and user data storage, MySQL serves as a complementary relational database.

2. Backend Services: The backend employs Spring Boot and Spring Data frameworks, a layered architecture with controllers, repositories, services, and transactional management for data processing and API communication.

3. Frontend Interface: The user interface, implemented in Angular, incorporates Angular Material and Tailwind CSS for styling and D3.js for interactive graph visualizations. Markdown and LaTeX support enhance the display of textual and mathematical content, enriching user experience.

4. Recommendation and Retrieval Layers: A hybrid approach combines Cypher-based graph queries and SQL analytics. Neo4j handles complex citation relationships, while MySQL provides statistical insights into auxiliary metadata.

5. Cloud Deployment: AWS services host the platform. Elastic Compute Cloud (EC2) manages backend services and AWS Lambda handling computationally intensive tasks like model inference. Neo4j runs on AuraDB, and Docker containerization ensures flexibility across different deployment environments.

B. Application Features and Results

The system’s effectiveness is showcased through several key functionalities, supported by visual demonstrations:

1) Vectorized Search: One of the core functionalities of our system is vectorized search. By entering a keyword like “GPT-3,” the system retrieves a list of relevant papers, including those directly focused on GPT-3 as well as related topics like large language models, prompt engineering, and parameter-efficient fine-tuning methods. Figure 4 illustrates an example of vectorized search in action.

2) Neo4j Visualization Preview: To enhance the understanding of relationships in the dataset, Neo4j visualizations provide an interactive view of the knowledge graph. Figure 5 shows a visualization where nodes represent papers, and edges indicate citation relationships.

3) Paper Example: BERT: The system provides detailed insights into individual papers. For example, as shown in Figure 6, the BERT paper display includes the abstract, author details, associated tasks (e.g., question answering, sentiment analysis), methods (e.g., residual connections, Adam optimizer), and links to related code implementations.

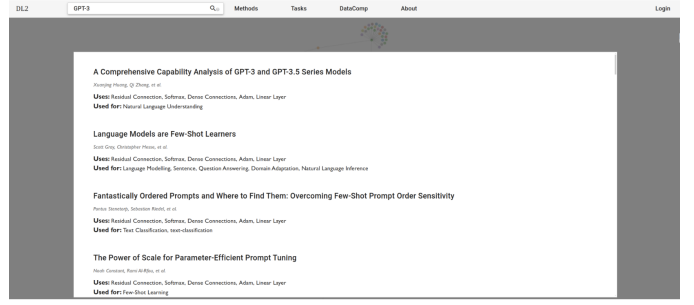


Fig. 4. An example of vectorized search for the keyword “GPT-3.”

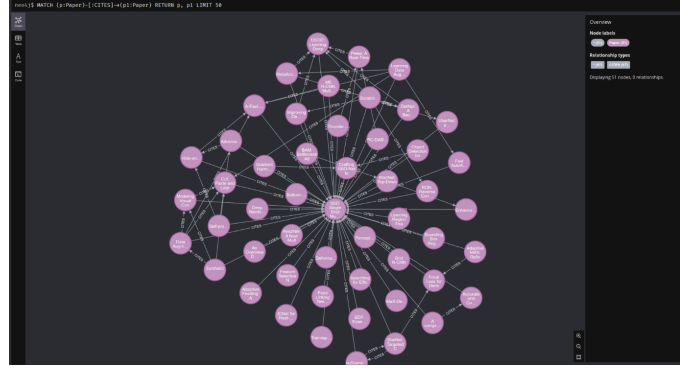


Fig. 5. Neo4j visualization preview of the knowledge graph. Nodes represent papers, and edges indicate citation relationships.

4) Open-ended Question Answering: The open-ended QA feature allows users to input queries in natural language and receive detailed responses. For example, Figure 7 illustrates a query about U-Net. It provides a comprehensive explanation of its architecture and applications.

5) Paper Summarization: The paper summarization feature extracts concise abstracts and highlights key details, such as core innovations, methodologies, and applications. As shown in Figure 8, this feature simplifies complex research, helping users quickly grasp a paper’s contributions.

C. Potential Bottlenecks and Improvements

Although the system’s architecture and functional design are sound, there are a number of possible bottlenecks and places that could use improvement. Below is a detailed discussion of these.

Neo4j’s graph query performance is a significant issue, especially when working with big subgraphs or sophisticated relationship searches. The computational cost of processing increases dramatically as the knowledge graph’s size and complexity rise. The user experience may suffer as a result of slower response times. In a similar vein, creating embeddings and retrieving vectors both need a significant amount of computing power. But as the dataset grows, it becomes more costly. This is particularly noticeable when dealing with large-scale retrieval tasks or high-dimensional embeddings.

Another issue with backend services is concurrency. The system might not be able to manage multiple data requests at

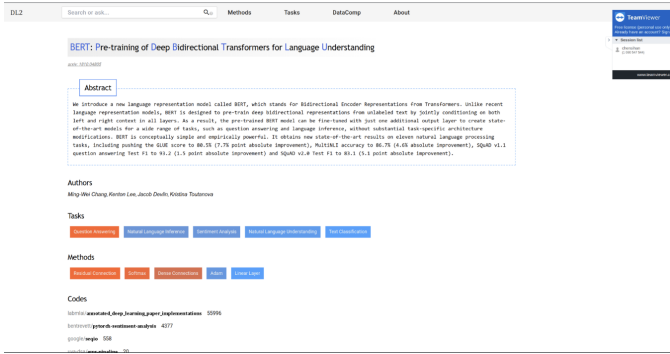


Fig. 6. An example of a paper display for the BERT model. It includes the abstract, associated tasks, methods, and related code implementations.

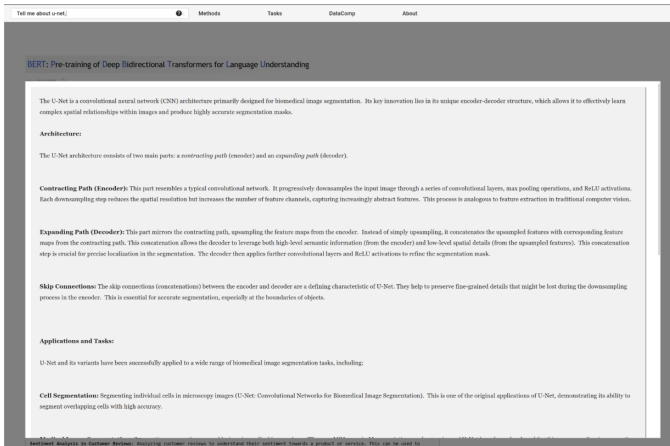


Fig. 7. An open-ended QA example where the system generates a detailed response about U-Net, including its architecture and applications.

once as the user base increases. Inadequate load control may cause service interruptions or delays for users.

Our system has undergone a number of enhancement techniques in order to overcome these issues. Neo4j’s indexing and caching methods can enhance query performance. The overhead of frequently retrieved data is decreased with this method. Efficiency could also be greatly increased by refining query pathways and using Neo4j’s integrated performance tweaking tools.

Using parallel indexing or switching to distributed vector databases helps reduce computing burden in embedding production and vector retrieval. These methods enable the system to handle bigger datasets and scale horizontally. The embeddings memories can be decreased by using sophisticated compression techniques.

Finally, load balancers and more sophisticated caching techniques can be used to improve backend services in order to handle growing concurrency. Incoming traffic can be equally distributed among servers by load balancers. This ensures that no one element is overpowered. Moreover, putting in place strong cache layers for static data helps lessen the strain on backend operations. A faster reaction can be generated as a result.

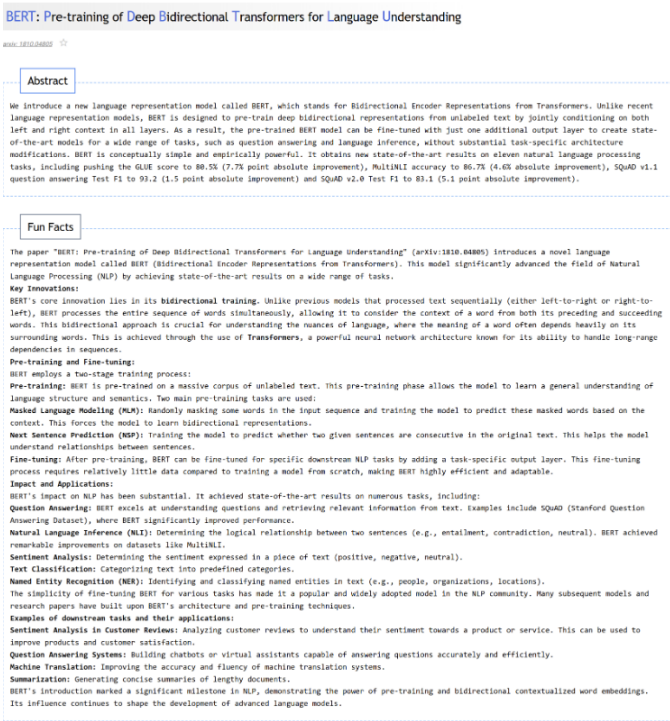


Fig. 8. An example of paper summarization, providing a concise abstract and key insights, such as methodologies and applications.

The system may maintain performance and scalability even as user demands and data complexity increase by methodically resolving these obstacles and putting the suggested enhancements into practice.

VI. EXPERIMENTS

The training process of the personalized learning platform’s graph attention network (GAT) model is pivotal in providing tailored recommendations for users. This section details the configuration employed during the training phase.

A. Training Setup

The training utilized a deep learning knowledge graph constructed from curated datasets. Node embeddings were precomputed with a dimensionality of 768 using external APIs. These embeddings served as direct inputs to the GAT model, bypassing the need for additional embedding layers.

B. Hyperparameters and Implementation Details

The GAT architecture included two layers, which is shown in Figure 9:

- The first layer transformed inputs from 768 to 1024 dimensions.
- The second layer projected 1024 dimensions back to 768 dimensions.

The model was trained for 100 epochs, utilizing a batch size of 128. The optimizer employed was Adam, with a learning rate set at 0.001. To balance computational efficiency and data diversity, the data loader restricted adjacency lists to

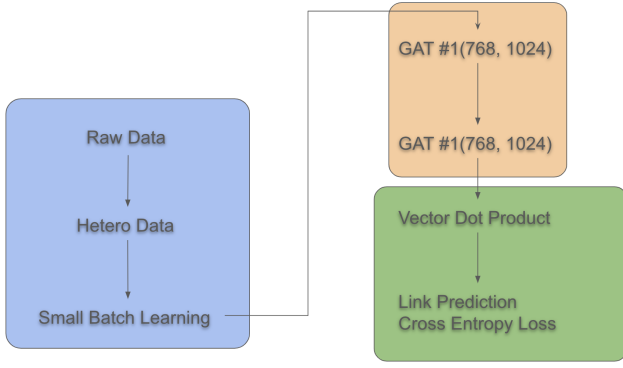


Fig. 9. Structure map for Graph Attention Network.

[50, 20, 10] neighbors for the first, second, and third layers, respectively.

C. Training Methodology

During training, the GAT model predicted user-paper interaction scores through dot products between user and paper embeddings. These scores were evaluated using binary cross-entropy loss, which incorporated a sigmoid function to normalize predictions. Early stopping was employed to prevent overfitting, halting training if validation loss did not improve for 10 consecutive epochs.

D. Results

The model converged stably, achieving a high accuracy on validation and test datasets. A visualization of loss trends across training iterations is depicted in Figure 10.

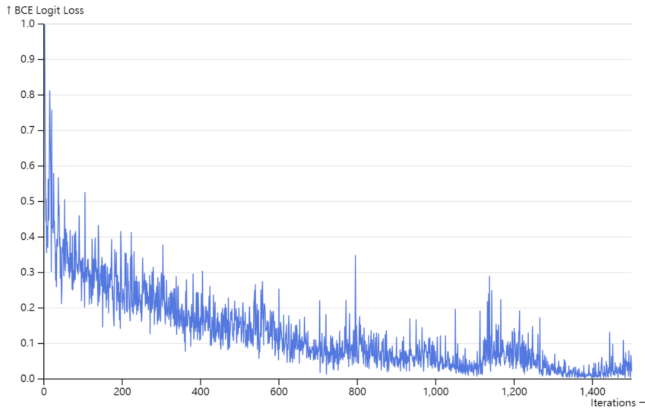


Fig. 10. Training and validation loss across 100 epochs.

The prediction task in this project focuses on edge classification, where the system identifies potential relationships between a user and candidate papers. To generate these candidate papers efficiently, the following heuristic subgraph search algorithm was employed:

- Compute the user's embedding vector based on previously liked papers.

- Retrieve papers with high similarity to the user's embedding.
- Expand the candidate set using the LinkNeighborLoader technique to query subgraphs within predefined constraints.

This process significantly narrows the search space for candidate papers and accelerates graph neural network inference.

To evaluate the model's predictive performance, the ROC (Receiver Operating Characteristic) curve and the AUC (Area Under the Curve) score were utilized. The ROC curve plots the True Positive Rate (TPR) against the False Positive Rate (FPR) across varying thresholds, defined as follows:

$$TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + TN},$$

where TP , FN , FP , and TN represent true positives, false negatives, false positives, and true negatives, respectively. The threshold function for classification is given by:

$$th(o, x) = \begin{cases} 1, & o > x \\ 0, & \text{otherwise.} \end{cases}$$

In the context of this project:

- **TP**: The model correctly predicts that a user likes a paper.
- **FN**: The model predicts that a user dislikes a paper, but the user actually likes it.
- **FP**: The model predicts that a user likes a paper, but the user actually dislikes it.
- **TN**: The model correctly predicts that a user dislikes a paper.

The AUC score represents the proportion of the area under the ROC curve. A random classifier achieves an AUC of 0.5, represented by a straight line connecting (0,0) and (1,1), while an ideal classifier achieves an AUC of 1. Figure 11 illustrates the ROC curve for the test set, with an AUC score of 0.91. The smooth curve indicates the model's ability to generalize well without overfitting.

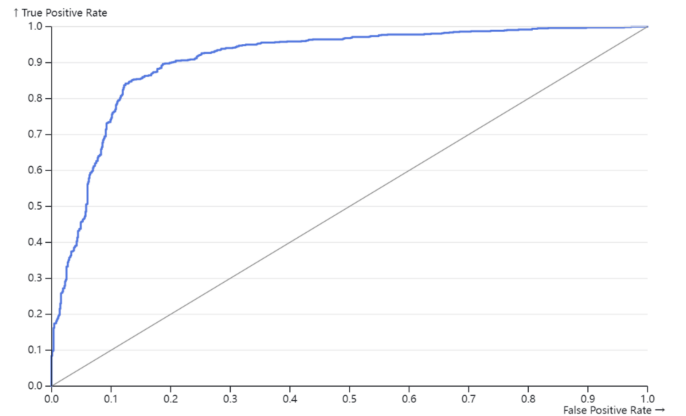


Fig. 11. ROC curve of the model's predictions on the test set, with an AUC score of 0.91.

E. Discussion

The experimental studies show how well the suggested two-layer Graph Attention Network (GAT) uses the deep learning knowledge graph and semantic embeddings for edge categorization. By forecasting user-paper interactions based on semantic and graph-based linkages, the model produces encouraging outcomes when offering research articles to users.

The model's extensibility and portability are two of its main advantages. Without retraining the entire model, the impacted subgraphs can be re-analyzed when new articles are uploaded to the system, updating the appropriate text embeddings and graph representations. In a similar vein, the system easily accepts new users. The model can produce user embeddings based on the semantic embeddings of research papers for tailored suggestions, regardless of the user's pre-existing preferences.

Furthermore, from a statistical and mathematical standpoint, the model connects "knowledge graphs" and "semantic embeddings" better. It supports the idea that a knowledge graph, which is a type of semantic web representation, captures both syntactic and semantic links between nodes. Semantic embeddings provide a strong framework for predictions and recommendations based on these graph interactions.

The method could be further enhanced in several areas despite its usefulness. First, when preferences change over time, adding user feedback into the training process may improve suggestion accuracy. Second, preserving scalability in practical applications would require optimizing the computational efficiency of extensive graph updates and embedding computations.

VII. CONCLUSION

In order to efficiently traverse and suggest scholarly research, this system created a personalized learning platform that integrates retrieval-augmented generation, graph neural networks, and a deep learning knowledge graph. The approach showed great promise in tackling the difficulties presented by dense citation networks and interdisciplinary research by combining semantic embeddings and graph-based linkages.

Promising outcomes came from the experiments, especially when it came to the Graph Attention Network (GAT), which predicted user-paper interactions with an AUC of 0.91. Additionally, by connecting retrieved papers to related themes and approaches, the system's integration with the knowledge graph improved interpretability, and it demonstrated its capacity to offer contextually appropriate recommendations. By providing thorough, succinct, and easily comprehensible insights into intricate research publications, it improved the user experience.

The value of integrating semantic embeddings with graph-based learning is among the most important lessons learned from this research. Furthermore, the system can easily handle additional users, articles, and changing research trends because to its modular architecture.

But the experiment also uncovered issues including embedding updates and computing limitations in large-scale graph queries. Limited GPU resources and complex front-end

programming to create an interference surface that is easy to use are additional challenges.

In order to further diversify the learning resources available to users, future work could investigate extending the platform to incorporate multimedia data, such as video lectures or presentations. Furthermore, adding multilingual features to the QA system can open up the platform to a wider audience. With these additions, the platform might become a comprehensive tool that can meet the various demands of both beginners and specialists.

REFERENCES

- [1] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [2] T. B. Brown, "Language models are few-shot learners," *arXiv preprint arXiv:2005.14165*, 2020.
- [3] E. Meijering, "A bird's-eye view of deep learning in bioimage analysis," *Computational and structural biotechnology journal*, vol. 18, pp. 2312–2325, 2020.
- [4] T.-Y. Ross and G. Dollár, "Focal loss for dense object detection," in *proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2980–2988.
- [5] W. Yu and X. Wang, "Mambaout: Do we really need mamba for vision?" *arXiv preprint arXiv:2405.07992*, 2024.
- [6] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.
- [7] E. W. Schneider, "Course modularization applied: The interface system and its implications for sequence control and data analysis." 1973.
- [8] A. Singhal, "Introducing the knowledge graph: things, not strings," Blog Post, Mountain View, CA, May 16 2012, <https://blog.google/products/search/introducing-knowledge-graph-things-not/>.
- [9] L. Ehlringer and W. Wöß, "Towards a definition of knowledge graphs," *SEMANTICS (Posters, Demos, SuCCESS)*, vol. 48, no. 1-4, p. 2, 2016.
- [10] A. McCallum *et al.*, "Cora dataset," Online, May 21 2017, <https://paperswithcode.com/dataset/cora>.
- [11] C. L. Giles *et al.*, "Citeseer dataset," Online, May 21 2017, <https://paperswithcode.com/dataset/citeseer>.
- [12] A. Farhangi *et al.*, "arxiv-10 dataset," Online, May 21 2022, <https://paperswithcode.com/dataset/arxiv>.
- [13] Z. Ye, Y. J. Kumar, G. O. Sing, F. Song, and J. Wang, "A comprehensive survey of graph neural networks for knowledge graphs," *IEEE Access*, vol. 10, pp. 75 729–75 741, 2022.
- [14] X. Wang, X. He, Y. Cao, M. Liu, and T.-S. Chua, "Kgat: Knowledge graph attention network for recommendation," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 950–958.
- [15] Y. Liu, S. Yang, Y. Xu, C. Miao, M. Wu, and J. Zhang, "Contextualized graph attention network for recommendation with item knowledge graph," *IEEE Transactions on knowledge and data engineering*, vol. 35, no. 1, pp. 181–195, 2021.
- [16] A. Vaswani, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017.
- [17] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [18] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand *et al.*, "Mixtral of experts," *arXiv preprint arXiv:2401.04088*, 2024.
- [19] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Akhtar, N. Barnes, and A. Mian, "A comprehensive overview of large language models," *arXiv preprint arXiv:2307.06435*, 2023.
- [20] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin *et al.*, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *ACM Transactions on Information Systems*, 2023.

- [21] H. Li, Y. Su, D. Cai, Y. Wang, and L. Liu, “A survey on retrieval-augmented text generation,” *arXiv preprint arXiv:2202.01110*, 2022.