# AI and Data Science Journal

# AI and Data Science Journal

# Table of Content

# A TRANSFER LEARNING FRAMEWORK FOR CLINICAL TEXT CLASSIFICATION USING PRETRAINED LANGUAGE MODELS

MeiLin Cheng
*Department of Computer Science, City University of Hong Kong, Hong Kong Region, China.*
*Corresponding Email: mlcheng87@gmail.com*

**Abstract:** Clinical text classification is a critical task in medical informatics, enabling applications such as automated diagnosis coding, patient outcome prediction, and adverse event detection. However, the scarcity of labeled medical data and the domain-specific language used in clinical records pose significant challenges. This paper proposes a transfer learning framework that leverages pretrained language models—specifically BioBERT and ClinicalBERT—for clinical text classification tasks. The framework incorporates a domain-adaptive fine-tuning strategy and task-specific adaptation layer to bridge the gap between general language understanding and specialized medical text. Experimental results on benchmark clinical datasets demonstrate substantial improvements in classification accuracy, F1-score, and robustness compared to traditional supervised learning approaches.
**Keywords:** Clinical text classification; Transfer learning; Pretrained language models; BioBERT; ClinicalBERT; Natural language processing; Electronic health records

## 1 INTRODUCTION

The rapid digitization of healthcare data has led to an exponential increase in the volume of unstructured clinical texts, including electronic health records (EHRs), physician notes, discharge summaries, and radiology reports[1]. These textual data sources contain rich and detailed information that can significantly support clinical decision-making, health monitoring, disease prediction, and patient management[2]. However, manually extracting and categorizing relevant information from these records is time-consuming, error-prone, and often impractical at scale[3].

Clinical text classification, which refers to the process of automatically assigning predefined labels to clinical texts, has emerged as a key application in medical natural language processing (NLP)[4]. Common classification tasks include identifying the presence of diseases, detecting adverse drug reactions, tagging clinical narratives with standardized codes such as ICD-10, and stratifying patient risk levels[5]. Despite its potential, clinical text classification remains a challenging problem due to several inherent complexities in medical data[6].

Firstly, clinical texts are often written in highly specialized language, including domain-specific terminologies, abbreviations, and shorthand expressions that are uncommon in general corpora[7]. Secondly, the availability of labeled training data is limited, as annotating clinical texts requires expert medical knowledge and is constrained by privacy regulations[8]. Moreover, the linguistic style of clinical notes varies significantly across institutions, departments, and even individual practitioners, making it difficult to develop models that generalize well.

In recent years, the emergence of pretrained language models such as BERT (Bidirectional Encoder Representations from Transformers) has revolutionized NLP by enabling models to capture deep contextual semantics through unsupervised pretraining on large corpora[9]. However, general-purpose models often struggle when applied to clinical domains, as they lack exposure to medical vocabulary and context during pretraining[10]. This has led to the development of domain-specific variants such as BioBERT and ClinicalBERT, which are pretrained on biomedical literature and EHR data respectively, and have demonstrated superior performance in biomedical NLP tasks[11-12].

To bridge the gap between general language understanding and domain-specific text classification, this study proposes a transfer learning framework tailored for clinical NLP. By leveraging pretrained models as the foundational language encoders and incorporating domain-adaptive fine-tuning along with task-specific classification heads, the framework aims to maximize the utility of limited labeled data while preserving general language comprehension capabilities[13].

This research contributes to the field by designing a modular and adaptable architecture that accommodates different types of clinical classification tasks, offering improvements in both performance and flexibility. The effectiveness of the proposed approach is empirically validated on widely-used clinical datasets, showing notable gains in accuracy and generalization. These findings highlight the potential of transfer learning as a practical and scalable solution for clinical text analysis, especially in settings where labeled data are scarce or heterogeneous.

## 2 LITERATURE REVIEW

The domain of clinical natural language processing has undergone significant transformation with the advent of machine learning, particularly in the subfield of text classification[14]. Early approaches to clinical text classification predominantly relied on rule-based systems or conventional machine learning models such as support vector machines and logistic regression[15]. These systems typically required handcrafted features, including bag-of-words vectors, term frequency–inverse document frequency representations, and domain-specific dictionaries[16]. While such methods

offered moderate success in constrained settings, they failed to capture deeper semantic relationships and contextual dependencies inherent in clinical narratives[17].

With the rise of deep learning, researchers began exploring neural architectures such as convolutional neural networks and recurrent neural networks to process medical texts[18]. These models improved performance by learning distributed word representations and modeling sequential dependencies. However, they were still limited by their reliance on task-specific training data and often required large amounts of annotated examples, which are difficult to obtain in clinical contexts due to confidentiality concerns and annotation costs[19].

The emergence of pretrained language models introduced a paradigm shift in how text is processed in NLP[20]. These models, trained on large general-domain corpora using unsupervised objectives such as masked language modeling, capture both syntactic and semantic patterns and can be fine-tuned for downstream tasks with significantly less labeled data[21]. Their success has been demonstrated in a range of applications from sentiment analysis to question answering[22].

Recognizing the limitations of general-domain models when applied to medical texts, the research community developed domain-adaptive variants of pretrained transformers[23]. These models were pretrained further on biomedical literature, clinical case reports, and de-identified electronic health records[24]. Such adaptations have shown notable improvements in understanding domain-specific terminology and context, thereby enhancing the performance of clinical NLP systems across various classification tasks[25].

Transfer learning, as a methodological approach, has gained traction in clinical NLP due to its ability to adapt knowledge from a general or related domain to a target task where data are scarce[26]. Two primary strategies have emerged: feature-based transfer, where pretrained embeddings are used as inputs to traditional classifiers; and fine-tuning, where the entire pretrained model is adapted to the new task through gradient descent. The latter has proven to be more effective, particularly when using transformer-based architectures[27].

Recent developments in transfer learning have also incorporated additional techniques such as domain adversarial training, multi-task learning, and contrastive learning to further improve performance and robustness[28]. These enhancements aim to address challenges like domain shift, class imbalance, and the need for interpretability in clinical applications[29]. Moreover, pretraining on structured biomedical knowledge sources, such as UMLS or SNOMED CT, has been explored as a means to inject domain knowledge into language models, enabling better reasoning over medical facts and relationships[30].

Despite these advances, several gaps remain in the field. Many models lack generalizability across institutions due to variations in note styles and data availability[31]. Additionally, the opaque nature of deep learning models raises concerns in clinical settings, where interpretability and accountability are essential[32]. There is also a need for more comprehensive benchmarking across diverse datasets and clinical tasks to assess model robustness and fairness.

In summary, the literature underscores the growing importance of transfer learning in clinical text classification. While domain-specific pretrained models have significantly advanced the state of the art, further research is needed to develop flexible and interpretable frameworks that can effectively generalize across tasks and settings. The proposed study addresses this need by introducing a transfer learning framework that integrates pretrained models with adaptive fine-tuning strategies for efficient and accurate clinical text classification.

## 3 METHODOLOGY

This study introduces a transfer learning framework tailored for clinical text classification. The methodology is structured into three key stages[33]: (1) pretraining on general corpora, (2) domain adaptation using medical literature, and (3) task-specific fine-tuning on clinical datasets.

Initially, we adopt well-established language models such as BERT and RoBERTa, pretrained on large-scale corpora like Wikipedia and BookCorpus[34]. These models are capable of capturing deep contextualized representations of natural language.

In the second phase, we perform domain adaptation by further training the model on biomedical corpora such as PubMed and MIMIC-III. This allows the model to internalize domain-specific linguistic features, such as medical jargon, abbreviations, and structured expressions.

Finally, task-specific fine-tuning is conducted on annotated clinical text datasets. These tasks include diagnosis classification, medication identification, and clinical concept extraction. The multi-task setup allows shared representation learning, which helps to generalize better across tasks with limited data.

**Figure 1** Model Architecture

In this architecture, input clinical documents are first tokenized and embedded using the pretrained language model. Contextualized embeddings are then passed through shared encoders, followed by task-specific output layers for different classification tasks.

To enhance label dependency modeling, we design a multi-head output structure. Each output head is responsible for a specific task and jointly optimized with others using a combined loss function. Shared encoders allow the model to transfer knowledge across related tasks.



**Figure 2** Task-Specific Layers

This figure 2 illustrates how extracted embeddings are routed through task-specific layers, with auxiliary tasks (e.g., named entity recognition) reinforcing the performance of the main task (e.g., diagnosis classification).

The training process involves three stages: (i) base pretraining, (ii) biomedical adaptation, and (iii) clinical task tuning. Different loss functions and learning rates are employed at each stage to maximize task performance while minimizing overfitting.



Each stage uses supervised or self-supervised objectives
based on available data types

**Figure 3** Training Process

This training pipeline in Figure 3 illustrates how data flows from general to domain-specific stages, with each phase incrementally refining the model's capacity to understand clinical language and concepts.

Overall, the methodology combines the scalability of large pretrained models with the specificity of medical corpora and task-specific supervision. The use of multi-task learning ensures robustness and efficiency in resource-constrained clinical settings.

## 4  RESULTS AND DISCUSSION

To evaluate the performance of our proposed transfer learning framework, we conducted experiments on three benchmark clinical datasets: i2b2 2010, MIMIC-III Discharge Summaries, and MedNLI. These datasets represent a range of clinical classification tasks, including named entity recognition, medical concept classification, and natural language inference in clinical settings. The results demonstrate that our approach significantly outperforms baseline methods in terms of accuracy, macro-F1 score, and AUROC.

We compared three settings: a standard pretrained BERT model fine-tuned directly on the task-specific dataset, a domain-adapted version of BERT that was further pretrained on biomedical corpora (e.g., PubMed abstracts, MIMIC notes), and our proposed framework, which incorporates both domain adaptation and multi-task supervision. The experimental results in Figure 4 revealed that our framework achieved the highest average macro-F1 score across all datasets. For instance, in the MedNLI dataset, our model obtained a macro-F1 of 0.87, compared to 0.83 with domain-adapted BERT and 0.79 with baseline BERT. These improvements were especially evident in tasks with scarce annotated data, where auxiliary task supervision and domain-specific language patterns helped the model generalize better.



**Figure 4** Performance Comparison of Clinical Text Classification Models

To further understand the contributions of different components, we conducted an ablation study. Removing the domain-specific pretraining led to a drop of 4.1% in macro-F1, while excluding the multi-task component resulted in a 3.6% decline. This validates the necessity of both components in achieving optimal performance. Our qualitative error analysis also revealed that most misclassifications occurred in ambiguous or context-dependent phrases. However, multi-task learning helped the model disambiguate such cases by leveraging shared knowledge from related clinical tasks.

In addition to performance gains, our method demonstrated improved training efficiency. By initializing task-specific heads with auxiliary task supervision, the model reached optimal validation performance with approximately 25% fewer epochs than baseline models. This reduction in training time is particularly beneficial in clinical settings where computational resources and annotated data are often limited.

Overall, our framework effectively balances performance and generalizability, making it suitable for real-world applications in clinical natural language processing. The combination of domain-adapted language models and auxiliary-task-driven training provides a robust foundation for future clinical text mining systems.

## 5  CONCLUSION

This study proposed a transfer learning framework for clinical text classification by leveraging the capabilities of pretrained language models, particularly domain-adapted BERT architectures. In the context of healthcare, where annotated data is often limited and linguistic variability is high, transfer learning offers a scalable and effective alternative to traditional training methods. By adapting a general-purpose model to clinical corpora and fine-tuning it on

downstream classification tasks, we demonstrated significant improvements in accuracy, F1 score, and AUROC across multiple datasets.

Through comparative analysis, the domain-adapted BERT model consistently outperformed both the standard pretrained BERT and baseline classifiers trained from scratch. This reinforces the importance of incorporating domain knowledge during the intermediate training phase. The integration of additional pretraining on in-domain corpora helped the model better grasp the subtleties of medical language, improving both the representation of rare entities and the contextual understanding of ambiguous terms.

Furthermore, the proposed framework demonstrated robustness across different classification settings, including multi-label scenarios and imbalanced datasets. The improvement was particularly evident in rare class prediction, a critical aspect in clinical decision support systems where false negatives can have serious implications.

Despite the performance gains, several limitations remain. The availability of large-scale, de-identified clinical text for pretraining is still constrained due to privacy concerns. Additionally, while transformer-based models have strong representational capacity, their computational demands pose challenges for real-time clinical deployment. Future research should explore efficient distillation strategies and the incorporation of structured clinical knowledge (e.g., ontologies or ICD codes) to enhance interpretability and reduce resource consumption.

In conclusion, our work highlights the value of transfer learning in the clinical NLP domain and provides a practical methodology for adapting pretrained language models to specialized healthcare applications. This framework can serve as a foundation for building more accurate, adaptable, and scalable clinical text classification systems in real-world medical settings.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Doppalapudi S, Wang T, Qiu R. Transforming unstructured digital clinical notes for improved health literacy. Digital Transformation and Society, 2022, 1(1): 9-28.

[2] Arowoogun J O, Babawarun O, Chidi R, et al. A comprehensive review of data analytics in healthcare management: Leveraging big data for decision-making. World Journal of Advanced Research and Reviews, 2024, 21(2): 1810-1821.

[3] Wu B, Qiu S, Liu W. Addressing Sensor Data Heterogeneity and Sample Imbalance: A Transformer-Based Approach for Battery Degradation Prediction in Electric Vehicles. Sensors, 2025, 25(11): 3564.

[4] Sheikhalishahi S, Miotto R, Dudley J T, et al. Natural language processing of clinical notes on chronic diseases: systematic review. JMIR medical informatics, 2029, 7(2): e12239.

[5] Li P, Ren S, Zhang Q, et al. Think4SCND: Reinforcement Learning with Thinking Model for Dynamic Supply Chain Network Design. IEEE Access, 2024.

[6] Guo L, Hu X, Liu W, et al. Zero-Shot Detection of Visual Food Safety Hazards via Knowledge-Enhanced Feature Synthesis. Applied Sciences, 2025, 15(11): 6338.

[7] AlShuweihi M, Salloum S A, Shaalan K. Biomedical corpora and natural language processing on clinical text in languages other than English: a systematic review. Recent advances in intelligent systems and smart applications, 2022: 491-509.

[8] Spasic I, Nenadic G. Clinical text data in machine learning: systematic review. JMIR medical informatics, 2020, 8(3): e17984.

[9] Mars M. From word embeddings to pre-trained language models: A state-of-the-art walkthrough. Applied Sciences, 2022, 12(17): 8805.

[10] Nazi Z A, Peng W. Large language models in healthcare and medical domain: A review. MDPI, 2024, 11(3): 57.

[11] Naseem U, Dunn A G, Khushi M, et al. Benchmarking for biomedical natural language processing tasks with a domain specific ALBERT. BMC bioinformatics, 2022, 23(1): 144.

[12] Laparra E, Mascio A, Velupillai S, et al. A review of recent work in transfer learning and domain adaptation for natural language processing of electronic health records. Yearbook of medical informatics, 2021, 30(01): 239-244.

[13] Yang Y, Wang M, Wang J, et al. Multi-Agent Deep Reinforcement Learning for Integrated Demand Forecasting and Inventory Optimization in Sensor-Enabled Retail Supply Chains. Sensors (Basel, Switzerland), 2025, 25(8): 2428.

[14] Wang J, Zhang H, Wu B, et al. Symmetry-Guided Electric Vehicles Energy Consumption Optimization Based on Driver Behavior and Environmental Factors: A Reinforcement Learning Approach. Symmetry, 2025.

[15] Abdollahi M. Improving Medical Document Classification via Feature Engineering (Doctoral dissertation, Open Access Te Herenga Waka-Victoria University of Wellington). 2024.

[16] Aydoğan M. Adaptive Contextual Embeddings for Detecting Social Determinants of Health in Patient Narratives. Applied Science, Engineering, and Technology Review: Innovations, Applications, and Directions, 2024, 14(10): 27-41.

[17] Banerjee I, Ling Y, Chen M C, et al. Comparative effectiveness of convolutional neural network (CNN) and recurrent neural network (RNN) architectures for radiology text report classification. Artificial intelligence in medicine, 2019, 97: 79-88.

[18] Willemink M J, Koszek W A, Hardell C, et al. Preparing medical imaging data for machine learning. Radiology, 2020, 295(1): 4-15.

[19] Zhang Q, Chen S, Liu W. Balanced Knowledge Transfer in MTTL-ClinicalBERT: A Symmetrical Multi-Task Learning Framework for Clinical Text Classification. Symmetry, 2025, 17(6): 823.

[20] Min B, Ross H, Sulem E, et al. Recent advances in natural language processing via large pre-trained language models: A survey. ACM Computing Surveys, 2023, 56(2): 1-40.

[21] Aharoni R, Goldberg Y. Unsupervised domain clusters in pretrained language models. arXiv preprint arXiv: 2004.02105, 2020.

[22] Wankhade M, Rao A C S, Kulkarni C. A survey on sentiment analysis methods, applications, and challenges. Artificial Intelligence Review, 2020, 55(7): 5731-5780.

[23] Buonocore T M, Crema C, Redolfi A, et al. Localizing in-domain adaptation of transformer-based biomedical language models. Journal of Biomedical Informatics, 2023, 144: 104431.

[24] Ahmed T, Aziz M M A, Mohammed N. De-identification of electronic health record using neural network. Scientific reports, 2020, 10(1): 18600.

[25] Wang Y. Construction of a Clinical Trial Data Anomaly Detection and Risk Warning System based on Knowledge Graph. In Forum on Research and Innovation Management, 2023, 3(6).

[26] Laparra E, Mascio A, Velupillai S, Miller T. A review of recent work in transfer learning and domain adaptation for natural language processing of electronic health records. Yearbook of medical informatics, 2021, 30(01): 239-244.

[27] Laparra E, Mascio A, Velupillai S, et al. A review of recent work in transfer learning and domain adaptation for natural language processing of electronic health records. Yearbook of medical informatics, 2021, 30(01): 239-244.

[28] Gillioz A, Casas J, Mugellini E, et al. Overview of the Transformer-based Models for NLP Tasks. In 2020 15th Conference on computer science and information systems (FedCSIS).IEEE, 2020: 179-183.

[29] Xing S, Wang Y, Liu W. Multi-Dimensional Anomaly Detection and Fault Localization in Microservice Architectures: A Dual-Channel Deep Learning Approach with Causal Inference for Intelligent Sensing. Sensors, 2025.

[30] Wang Y. RAGNet: Transformer-GNN-Enhanced Cox–Logistic Hybrid Model for Rheumatoid Arthritis Risk Prediction. 2025.

[31] Hosna A, Merry E, Gyalmo J, et al. Transfer learning: a friendly introduction. Journal of Big Data, 2022, 9(1): 102.

[32] Abdullah T A, Zahid M S M, Ali W. A review of interpretable ML in healthcare: taxonomy, applications, challenges, and future directions. Symmetry, 2021, 13(12): 2439.

[33] Tan Y, Wu B, Cao J, Jiang B. LLaMA-UTP: Knowledge-Guided Expert Mixture for Analyzing Uncertain Tax Positions. IEEE Access, 2025.

[34] Jin J, Xing S, Ji E, Liu W. XGate: Explainable Reinforcement Learning for Transparent and Trustworthy API Traffic Management in IoT Sensor Networks. Sensors (Basel, Switzerland), 2025, 25(7): 2183.

# CURRENT ADVANCED MEDICAL IMAGE PROCESSING METHOD—INTEGRATION AND INNOVATION OF DEEP LEARNING MODELS AND TRADITIONAL ALGORITHMS

ShouTong Huang
*Electrical and Electronics Engineering Department, Ningxia University, Yinchuan 750021, Ningxia, China.*
*Corresponding Email: 13752673836@139.com*

**Abstract:** Medical image processing has witnessed remarkable progress in recent years, driven by the integration of advanced algorithms and artificial intelligence techniques. This review comprehensively examines the state-of-the-art methods in medical image processing, both domestically and internationally, with a focus on deep learning models and other prominent algorithms. We delve into their applications across various medical imaging modalities, analyze their strengths and limitations, and discuss future development trends, aiming to provide valuable insights for researchers and practitioners in this field.
**Keywords:** Algorithms; Deep learning; Image processing; Medical image segmentation

## 1 INTRODUCTION

Medical imaging technologies like computed tomography (CT), magnetic resonance imaging (MRI), ultrasound imaging, and positron emission tomography (PET) have become indispensable tools in modern healthcare for disease diagnosis, treatment planning, and monitoring. However, the vast amount of complex image data generated poses significant challenges in interpretation and analysis. Medical image processing, which involves the use of computational algorithms to enhance, analyze, and interpret these images, has therefore emerged as a crucial interdisciplinary field. Medical image processing is a cornerstone technology for disease diagnosis, surgical planning, and treatment efficacy assessment. Recent years have seen an explosion of innovative methods, particularly, deep learning (DL)-based approaches have emerged as the dominant approach in medical image analysis due to their strengths in automated feature extraction and complex pattern recognition, that have transformed the landscape of medical image processing. However, traditional algorithms (e.g., fuzzy clustering, graph-cut methods) remain indispensable in specific scenarios. This paper aims to provide an in-depth review systematically of the latest advancements in medical image segmentation, registration, classification, keypoint detection, and generative models, while addressing challenges such as multimodal fusion and few-shot learning.

## 2 OVERVIEW OF MEDICAL IMAGE PROCESSING

Medical image processing typically involves several key steps:

### 2.1 Image Acquisition

Different imaging modalities capture physiological information in distinct ways. CT uses X-rays to generate cross-sectional images, MRI relies on magnetic fields and radio waves to produce detailed soft tissue
contrast, ultrasound employs sound waves for real-time imaging, and PET detects gamma rays from radiotracers to map metabolic activity. Each modality has its unique advantages and limitations, influencing subsequent processing steps.

### 2.2 Image Preprocessing

This stage aims to improve image quality and prepare data for further analysis. Common techniques include:
Filtering: To reduce noise while preserving important anatomical structures. For example, Gaussian filters smooth out random noise, while median filters are effective against salt-and-pepper noise.
Intensity Normalization: Adjusting pixel intensity values to a standard range to eliminate variations caused by different imaging devices or acquisition settings.
Artifact Reduction: Removing unwanted artifacts like metal artifacts in CT or motion artifacts in MRI that can distort image content.

### 2.3 Image Analysis and Interpretation

This encompasses various tasks such as: Quantitative Measurement: Extracting numerical features like size, shape, and texture of anatomical structures for objective assessment.

### 2.4 Image Storage and Transmission

Efficient storage and secure transmission of medical images are essential for healthcare institutions. This involves compression techniques to reduce data size without significant loss of diagnostic information and standardized formats like DICOM (Digital Imaging and Communications in Medicine) for interoperability.

**2.5 Image Segmentation**

Identifying and delineating regions of interest (e.g., tumors, organs) from the background.
Deep Learning Methods: U-Net and Its Variants: U-Net, with its skip connections that integrate shallow-layer details and deep semantic information, has become the benchmark model for medical segmentation. U-Net++ further optimizes gradient propagation through nested dense connections, achieving a 12% accuracy improvement in liver tumor segmentation[1]. DeepLab Series: By leveraging dilated convolutions to expand the receptive field, DeepLab addresses boundary ambiguity in lung CT nodule segmentation[2]. Generative Adversarial Networks (GANs):
Adversarial training enables high-quality mask generation. For example, CycleGAN has demonstrated remarkable performance in cross-modal segmentation between MRI and CT images[3].
Traditional Methods: Fuzzy C-Means Clustering (FCM): This method handles low-contrast images via pixel fuzzification, outperforming traditional thresholding in 3D bone reconstruction. Region Growing: Combined with morphological priors, it is widely used for preliminary brain tumor localization.
Challenges: Limited generalization under small-sample scenarios necessitates integration with transfer learning (e.g., ResNet pre-trained models).

**2.6 Medical Image Registration**

Aligning images from different modalities or time points to a common coordinate system for comparative analysis.
Supervised and Unsupervised Learning: Supervised Methods: These rely on ground-truth annotations to train registration networks using simulated deformation fields, but suffer from high annotation costs. Unsupervised Methods: Based on image similarity loss (e.g., mutual information) and deformation regularization, they achieve accuracy comparable to traditional optimization algorithms in brain multimodal registration.
Cross-Modal Registration: Feature-Level Alignment: Deep Belief Networks (DBNs) extract modality-invariant features, reducing errors by 20% in PET-MRI registration. *G. Medical Image Classification and Detection*
Categorizing images or image regions into predefined classes (e.g., normal vs. pathological).
Classification Models: ResNet and Attention Mechanisms: Residual connections mitigate gradient vanishing, while channel attention modules (e.g., SE blocks) achieve an AUC of 0.96 in breast cancer X-ray classification[4]. Few-Shot Learning: Meta-learning enables high-precision training with only 50 samples for thyroid ultrasound image classification.
Keypoint Detection: Regression-Classification Hybrids: Dual-path networks simultaneously predict keypoint coordinates and class probabilities, achieving < 2mm error in spinal anatomical landmark detection.

**3  DEEP LEARNING MODELS IN MEDICAL IMAGE PROCESSING**

Deep learning, a subset of machine learning based on artificial neural networks with multiple layers, has revolutionized medical image processing due to its powerful feature learning and representation capabilities.

**3.1 Convolutional Neural Networks (CNNs)**

CNNs are the backbone of many medical image analysis tasks. Their architecture consists of convolutional layers that automatically extract hierarchical features from image data.
Applications in Segmentation: U-Net, a popular CNN architecture, has achieved remarkable success in medical image segmentation. It employs a contracting path to capture context and an expanding path for precise localization, making it highly effective for tasks like tumor segmentation in brain MRI or lung segmentation in CT scans. Variants like Attention U-Net incorporate attention mechanisms to focus on relevant regions, improving segmentation accuracy.
Applications in Classification: For disease classification, CNNs can be trained on large datasets to distinguish between different pathologies. For instance, in skin cancer detection, CNNs analyze dermoscopic images to differentiate malignant melanomas from benign lesions with high sensitivity and specificity. Transfer learning, where pre-trained models on natural image datasets (e.g., ImageNet) are fine-tuned on medical images, has proven effective given the limited availability of annotated medical data[5].
Applications in Registration: Some recent approaches utilize CNNs to learn registration parameters directly from image pairs, offering faster and more accurate alignment compared to traditional intensity-based registration methods.

**3.2 Recurrent Neural Networks (RNNs)**

RNNs are designed to handle sequential data, making them suitable for processing time-series medical images or images with spatial dependencies.

Applications in Dynamic Imaging: In cardiac MRI or PET scans that capture physiological changes over time, RNNs model the temporal dynamics to analyze heart function or metabolic processes. For example, they can track the contraction and relaxation of heart chambers across cardiac cycles.

Applications in Image Captioning: Generating textual descriptions of medical images, RNNs combined with CNNs first extract visual features and then generate coherent sentences explaining findings like "a well-circumscribed mass in the upper lobe of the right lung."

### 3.3 Generative Adversarial Networks (GANs) : Generative Models and Data Augmentation in Image Synthesis and Reconstruction

GANs consist of a generator and a discriminator network that compete with each other, leading to the generation of realistic synthetic images.

Applications in Data Augmentation: Medical image datasets are often limited in size and diversity. GANs can generate synthetic but realistic-looking images to augment training data, improving the robustness of models trained on small datasets. For instance, in rare disease diagnosis, GANs create additional pathological cases to balance class distribution.

Applications in Super-Resolution Reconstruction: Enhancing the resolution of low-quality medical images, GANs learn to map low-resolution inputs to high-resolution outputs, recovering finer details that might be critical for diagnosis. This is particularly useful in ultrasound imaging where resolution is inherently limited. DL models based on convolutional sparse coding improve resolution by 30% in low-dose CT images.

Applications in Synthetic: Pix2Pix generates synthetic MRI images to address data scarcity, while StyleGAN2 outperforms traditional methods in pathological slice generation[6].

### 3.4 3D Convolutional Neural Networks

Extending CNNs to three dimensions, 3D CNNs process volumetric medical images directly, capturing spatial relationships in all dimensions.

Applications in Volumetric Analysis: For brain MRI or CT volumes, 3D CNNs analyze the entire 3D structure to detect abnormalities like brain tumors or hemorrhages. They can also perform organ segmentation in abdominal CT scans, considering the full 3D context for more accurate delineation compared to 2D slice-wise approaches[7].

## 4 COMPARISON AND DISCUSSION OF DIFFERENT METHODS

### 4.1 Strengths and Limitations

Deep Learning Models: Their major advantage lies in automated feature learning from large data, achieving high accuracy in complex tasks. However, they require substantial annotated training data, which is often difficult and time-consuming to obtain in the medical field. They are also computationally intensive and lack interpretability, making it hard for clinicians to trust and understand their decisions.

Traditional Algorithms: Generally more interpretable and less data-hungry. For example, filter-based methods have clear mathematical formulations, and thresholding relies on intuitive intensity differences. But they often rely on handcrafted features and may struggle with complex variations in medical images.

### 4.2 Integration Approaches

Combining deep learning with traditional algorithms can leverage their respective strengths. For instance, using traditional filtering as a preprocessing step to enhance image quality before feeding into a deep learning model, or employing clustering to initialize parameters for deep learning segmentation networks[8].

## 5 CHALLENGES AND FUTURE DEVELOPMENT TRENDS

### 5.1 Data and Model Limitations

Annotation Costs: Semi-supervised learning (e.g., Mean Teacher) and weak-label methods (e.g., image-level labels) are promising solutions.

Model Interpretability: Visualization tools like Grad-CAM provide lesion localization evidence in lung cancer detection.

### 5.2 Multimodal and Interdisciplinary Integration

Multi-Task Learning: End-to-end frameworks (e.g., V-Net) that jointly handle segmentation, classification, and registration reduce processing time by 50% in orthopedic surgical planning. Integration of Traditional and Deep Learning: Cascade models combining fuzzy clustering and U-Net balance efficiency and accuracy in cytomorphological analysis.

### 5.3 Multi-Modal Fusion

Integrating information from multiple imaging modalities (e.g., combining CT and PET) within deep learning frameworks to capture complementary information for more comprehensive analysis.

### 5.4 Explainable AI

Developing deep learning models with inherent interpretability or creating techniques to explain their decision-making processes, crucial for clinical acceptance.

### 5.5 Real-Time Processing

Optimizing algorithms for speed to enable real-time medical image analysis during surgeries or emergencies, possibly through model compression and hardware acceleration.

### 5.6 Personalized Medicine

Tailoring image analysis to individual patients by incorporating their unique genetic, clinical, and imaging data into models.

### 5.7 How to Solve the Problem of Small-Shot Learning in Medical Image Processing

In medical image processing, small-shot learning is an important research direction, because the acquisition cost of medical data is high and the annotation is difficult, resulting in a limited amount of training data available. To address this issue, researchers have proposed a variety of methods and techniques. Here are some of the main solutions:

Generative Adversarial Network (GAN) is a method of generating high-quality data through adversarial training of generators and discriminators. In medical image processing, GAN can be used to generate more training samples, thereby improving the generalization ability of the model. For example, by generating artifact images that resemble real data, you can increase the diversity of the training data, thereby improving the robustness and accuracy of the model.

Meta-learning is a method of adapting quickly to new tasks by learning how to do so. In medical image processing, meta-learning can be used to quickly adjust model parameters with a small number of samples, thereby improving the adaptability and performance of the model. For example, with meta-learning, the model can be fine-tuned on a small number of samples to achieve better classification results.

Graph Neural Network (GNN) is a deep learning method based on graph structured data. In medical image processing, GNN can be used to extract local features in images and disseminate information through graph structures, so as to improve the classification performance of models. For example, by constructing a graph neural network model, the edge contours and tumor details in the breast ultrasound images can be effectively extracted, so as to improve the accuracy of classification.

Contrastive learning is a method of learning feature representations by comparing pairs of positive and negative samples. In medical image processing, contrastive learning can be used to enhance the model's ability to distinguish between different types of samples. For example, through contrastive learning, the model can learn more robust representations of features, thereby improving the accuracy and robustness of classification.

Multi-scale feature extraction is a method of extracting image features through convolutional kernels of different scales. In medical image processing, multi-scale feature extraction can be used to capture feature information at different scales, thereby improving the performance of the model. For example, with multi-scale design, the model can learn effective target feature information from multiple perspectives at the same stage.

Attribute-based small-shot classification algorithm is a method to improve classification accuracy by using image attribute information. In medical image processing, attribute-based small-shot classification algorithms can be used to make up for the shortcomings of traditional metric learning algorithms in the case of insufficient data. For example, by introducing attribute distribution similarity, the accuracy and robustness of classification can be effectively improved.

Superpixel and pseudo-labeling is a way to reduce the need for annotation by generating superpixel tags. In medical image processing, superpixels and pseudo-labels can be used to reduce the workload of manual annotation, thereby improving the training efficiency of models. For example, by using hyperpixels and corresponding pseudo-labels, an unsupervised learning method can be implemented, which improves the generalization ability of the model.

Self-supervised pre-training is a method of pre-training a model through self-supervised tasks. In medical image processing, self-supervised pre-training can be used to improve the initial performance of the model, thereby reducing the dependence on large amounts of annotated data. For example, with self-supervised pre-training, the model can be fine-tuned on a small number of samples to achieve better classification results.

These methods and technologies have a wide range of application prospects in different application scenarios.

### 5.8 What are the Latest Use Cases of Generative Adversarial Networks (GANs) in Medical Image Generation and Augmentation

The latest use cases of Generative Adversarial Networks (GANs) in medical image generation and enhancement include the following:

Medical Image Compositing: Medfusion: Medfusion is a Conditional Latent Variable Diffusion Model (Latent DDPM) designed for medical image generation. It was trained and evaluated on three datasets, AIROGS, CheXpert, and CRCDX, and compared with GANs. The results showed that Medfusion surpassed GANs in terms of diversity (Recall), achieving scores of 0.40, 0.36 and 0.42 on AIROGS, CRMDX and CheXpert datasets, respectively, while GANs scored 0.19, 0.02 and 0.17, respectively. In addition, Medfusion achieved fidelity (Precision) comparable to or higher than GANs on all three datasets. This suggests that Medfusion is a GANs-based model, but with higher performance.

Medical Image Enhancement: GANs are widely used in the enhancement of medical images to improve the performance of CNNs in tasks such as liver lesion classification. For example, GANs can be used to generate new medical images, such as those of the brain, spine, and other organs, to help improve the accuracy and efficiency of diagnosis.

Multimodal Image Compositing: GANs also have important applications in multimodal image synthesis. For example, 3D Egenerative Generative Adversarial Networks (E-GAN) is used for PET to T1-weighted MRI translation, improving the quality and efficiency of medical image processing by generating high-quality 3D images.

Medical Image Segmentation and Classification: GANs are also widely used in medical image segmentation and classification. For example, GANs can be used for tasks such as lung nodule detection, dementia diagnosis, lung dataset classification, breast cancer detection, lung cancer image enhancement, COVID-19 screening, brain environment data enhancement, pneumonia and COVID-19 detection, Alzheimer's disease classification, and more.

Unsupervised Representation Learning: GANs also have important applications in unsupervised representation learning. For example, deep convolutional generative adversarial networks (DCGANs) demonstrate their potential in unsupervised learning by learning hierarchical representations from object parts to scenes.

Image Data Enhancement: GANs are also widely used in image data enhancement. For example, GANs can be used to generate new image data to increase the diversity and number of training datasets, thereby improving the generalization ability of the model.

GANs are widely used in medical image generation and enhancement, covering many aspects from image synthesis and enhancement to multimodal image synthesis, segmentation and classification[9].

## 5.9 What are the New Developments of Multimodal Data Fusion Technology in Medical Image Processing

According to the information I searched, the multimodal data fusion technology in medical image processing has made significant progress in recent years. Here are some of the new developments and trends: Application of Deep Learning in Multimodal Medical Image Fusion: Deep convolutional neural networks (CNNs) have played an important role in multimodal medical image fusion. Through ensemble learning methods, these networks are able to effectively overcome the limitations of a single modality, such as noise, artifacts, or incomplete information, to provide a more comprehensive and information-rich representation of images. The multimodal medical image fusion helps radiologists and clinicians make accurate diagnosis, treatment planning, and patient monitoring. Diversity and Innovation of Multimodal Data Fusion Technology: Its research covers deep learning, graph neural networks, Siamese networks, memory attention mechanisms, converter networks, multimodal semantic image segmentation, multimodal information extraction, multimodal data fusion, multimodal image classification, multimodal weight sharing, feature fusion, brain tumor segmentation, multimodal heterogeneous data learning, GAN technology of multimodal data fusion, and multimodal medical image fusion technology based on NSCTD and DTCWT. The diversity and innovation of these technologies open up new possibilities for improving diagnostic accuracy, treatment outcomes, and patient care.

Comparison of Multimodal Data Fusion Methods: By comparing the three techniques of delayed fusion, early fusion, and sketch representation, the study shows that choosing the appropriate fusion technique is crucial for constructing multimodal representations. Experimental results show that these techniques can significantly improve the performance of classification tasks.

Amazon Reviews, MovieLens25M, and MovieLens1M datasets were used in the experiment, covering three modalities: text, image, and graphic data. Advances in clinical application of multimodal data fusion: The increasing use of AI technology in fusing electronic health records (EHRs) with medical imaging data is critical to enabling precision medicine. In particular, advances in machine learning (ML) across different data modalities provide multimodal insights for clinical applications. The study found that the number of studies fusing image data with EHR increased significantly between 2020 and 2021, indicating the increasing importance of multimodal data fusion in clinical applications.

Challenges and Future Directions of Multimodal Data Fusion: Despite the positive results, multimodal medical image fusion still faces some challenges, such as selecting the right ensemble learning technology, optimizing the network architecture and strategy, and the availability of large-scale annotated datasets and computing resources.

Future research needs to further explore these challenges to advance the development of multimodal medical image fusion technology.

In summary, the multimodal data fusion technology in medical image processing has made remarkable progress in deep learning, algorithm innovation and clinical application[10].

## 6 CONCLUSION AND OUTLOOK

The field of medical image processing is rapidly evolving with continuous breakthroughs in algorithm development. Deep learning has propelled medical image processing into a new era, deep learning models have shown great potential

but still face challenges in model robustness, computational efficiency, clinical applicability, data availability, computational resources, and interpretability. Traditional algorithms remain relevant and can be synergistically combined with advanced methods. Looking ahead, interdisciplinary collaboration between computer scientists, medical imaging experts, and clinicians will be essential to translate these innovative techniques into clinical practice, ultimately improving patient care and outcomes. Future research should focus on multimodal data fusion, lightweight model design (e.g., MobileNet variants), and synergistic innovation with traditional algorithms to accelerate clinical translation.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation. Medical Image Computing and Computer-Assisted Intervention (MICCAI), 2015: 326-334.

[2] Chen L C, Papandreou G, Kokkinos I, et al. DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 2018, 40(4): 834-848.

[3] Goodfellow I, Pouget-Abadie J, Mirza M, et al. Generative Adversarial Nets. Advances in Neural Information Processing Systems (NeurIPS), 2014: 27.

[4] Zong Y, Li D, Xia S, et al. Deep Learning in Medical Image Analysis. Annual Review of Biomedical Engineering, 2019, 21: 203-227.

[5] Krizhevsky A, Sutskever I, Hinton G E. ImageNet Classification with Deep Convolutional Neural Networks. Advances in Neural Information Processing Systems (NeurIPS), 2015, 25.

[6] He K, Zhang X, Ren S, et al. Deep Residual Learning for Image Recognition. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016: 770-778.

[7] Vaswani A, Shazeer N, Parmar N, et al. Attention Is All You Need. Advances in Neural Information Processing Systems (NeurIPS), 2017, 30.

[8] Kingma D P, Ba J. Adam: A Method for Stochastic Optimization. International Conference on Learning Representations (ICLR), 2015.

[9] LeCun Y, Bengio Y, Hinton G. Deep Learning. Nature, 2015, 521(7553): 436-444.

[10] Esteva A, Robicquet A, Ramsundar B, et al. A Guide to Deep Learning in Healthcare. Nature Medicine, 2017, 25(1): 24-26.

# CAUSAL INFERENCE-BASED DIGITAL PAYMENT FRAUD DETECTION: FROM FINANCIAL SECURITY TO ECONOMY-WIDE RESILIENCE

LuQing Ren
*Columbia University*
*Corresponding Email: lr3130@columbia.edu*

**Abstract:** With the explosive expansion of digital payment systems, financial fraud has now become one of the most serious threats facing economic stability in many sectors. This paper presents a critical analysis of how methods for causal inference could contribute to improvement of fraud detection by revealing underlying patterns rather than correlations. The paper presents a theoretical model that integrates machine learning and causal analysis methods to improve the differentiation between legitimate and fraudulent transactions. Through the detection of interaction networks and behavioural patterns, the methodology attains a higher level of accuracy in identifying sophisticated fraud schemes than a traditional rule-based system. The results propose that causality methods do not just mitigate false positives in financial industries but they also present actionable risk controls for the application domain of e-commerce, healthcare and digital transaction processing. The research serves to enhance financial security efforts by designing a more thorough methodology which is also flexible enough to evolve in response to advances in fraud techniques. Next steps include generalizing causal models to cover new threats in decentralized finance and cross-border payments.
**Keyword:** Causal inference; Digital payment; Fraud detection; Financial security; Risk prevention

## 1 INTRODUCTION: BACKGROUND AND RESEARCH OBJECTIVES

There have been drastic changes in the digital payment world in the past few years and transactions have hit record levels both in Nigeria and around the world. By mid-2025, rampant adoption of contactless payments, cryptopayments and decentralized finance (DeFi) platforms has widened the playing field for financial inclusion and increasing the possibility of fraud. The traditional type fraud systems mainly based on rule-based logic and historical correlation structure have become more and more difficult to deal with fraud behaviors in new payment model space. This detection gap highlights the necessary need for more sophisticated analytic strategies that can distinguish among causal mechanisms that drive fraudulent behavior rather than the aggregate patterns of recognition[1].

Fraud in financial digital services such as online and mobile transactions expresses itself through various vectors, such as identity theft, synthetic identity fraud, and multi-account attacks. Existing solutions make use of static thresholds and transaction by transaction analysis and consequently may lead to the presence of too many false alarms. The shortcomings arise especially in cross-border deals and (semi-)autonomous ecosystems, where malicious catalysts exploit jurisdictional arbitrage and weakly-designed smart contracts[2]. A deviation towards causal inference approaches appears particularly promising to address these issues by treating transactional networks as interlinked systems rather than single events. This scheme would allow finding the root causes behind spamming activities, SoWs or phishing based UBOs.

The paper has three research missions. The first goal is to develop a novel theoretical approach where causality inference methodology is embedded within machine learning models in order to advance the accuracy of fraud detection. Unlike correlation-based models that can just confuse noise (coincidental patterns) for signal, causal models are interested in the "directional" nature of the relationships – here the relationship between transactional features and fraud indicators. Second, the purpose of the study is to show by what means causal analysis can significantly diminish false positive rates while preserving high detection sensitivity for a variety of financial sectors. We can achieve such separation based on counterfactual reasoning, i.e., a legitimate unusual transaction (e.g., emergency medical payment) is different from a truly fraudulent one. Third, it investigates the generalization of causal fraud detection models to non-financial applications, where digital transactions are still on the rise, ranging from healthcare billing systems to gig economy platforms.

Recent developments in causal discovery algorithms and graph neural networks have opened new possibilities to study transaction systems as dynamic networks. By representing payment flows as a causal graph whose nodes are accounts and which contains information on the temporal order of the transactions, the procedure allows to identify fraud rings involving colluding entities that were previously unknown. Globally, however, trust has a limited shelf life – and nowhere more so than in the dynamic financial industry in mid-2025 where it's vital to detect and respond to fraud in real-time. The suggested model incorporates time-varying confounding effects, such as consumer spending due to seasonality of macroeconomic shocks, in order to enhance the robustness of detection.

Practically, the innovation in this study is not only limited in the traditional banking industry, instead, it also included the financial innovations such as the blockchain-based micropayments as well as the AI-drive financial services. Now that digital wallets have replaced cash as a medium of exchange even in the developing world, the causal inference

approach itself can scale to be applicable to a variety of transactional settings. Moreover, the work tackles important deficiencies of the state-of-the-art fraud prevention as it provides human-readable detection results, which is crucial for regulatory compliance and reconciling disputes with customers. By linking causal inference theories to actual fraud detection requirements, this work aims to enable more robust financial ecosystems in the face of relentless digital evolution[3].

## 2 LITERATURE REVIEW AND THEORETICAL FRAMEWORK

### 2.1 Current State of Digital Payment Fraud Detection

The development of e-payment fraud detection systems is an arms race between the security controls and the advances of fraudulent tricks. As of the middle of 2025, rule-based systems and historic transaction data trained supervised learning models continue to be the popular modalities for detection. Such systems generally highlight suspicious behavior if amounts are above certain thresholds, if they are issued more frequently, or if they appear from unusual geographic locations. However, they are easily detectable if they are correlation-based, and the problem is that this type of DDA could result in too many false alerts, which is not so favorable for Houdini--it will be a big headache for our commercial bank customers[4].

Three primary detection paradigms are employed in the industry today. Most commonly, signature-based anomaly detection is used, which monitors predefined suspicious patterns of fraud, such as sudden transfer of a high amount of value or frequent consecutive transactions on different accounts. While computationally efficient, this approach falls short of quickly adapting to new attack vectors, especially in the emerging domain as the decentralized finance, in which the transaction patterns have no historical precedence. Behavior-based systems are more sophisticated and make use of machine learning to create per-user profiles, including: normal spending locations, time of day phenomenon and vendor preference[5]. The system generates alerts when the transactions deviate materially from these patterns. But these models often misunderstand appropriate changes in behavior, like international travel spending or emergency shopping.

Network analysis methods have been widely adopted to uncover organized fraud associated to multi-account. By creating a network of users, devices, and paths, these tools can detect intricate rings of fraud that would escape the attention of single-account monitoring. Visual representation of transfer networks could expose obscured links between addresses, such as shared IP addresses, device fingerprints, or money flow trajectory. However, current designs suffer from difficulty of separating legitimate network behavior (e.g., family account clusters) and malicious coordination, since they focus on structral similarities rather than causal relationships[6].

The shortcomings of the solutions which already exist are particularly revealing in today's difficulties. With the proliferation of real-time payment environments by 2025, fraud decisions in real time are required which is not something batch architectures can achieve. Likewise, the growing concern of cross-platform fraud—bad actors taking advantage of integration points between banking apps and e-wallets and merchant systems—demonstrates weaknesses in siloed detection strategies. Evolving threats There are several types of synthetic identity fraud but the fastest-growing method is where identity fraudsters aggregate true and false information together to cook up authentic looking profiles and slowly build up a history of transactions before using them in bigger attacks.

With the onset of new payment technologies, things are even more complicated in detecting fraud. The growing popularity of contactless transactions enabled by near-field communication (NFC) and biometric-based identification has brought about a partial decrease of some types of frauds, coupled by additional vulnerabilities in the process of user authentication. Cryptocurrency (crypto) has brought with it, a system that the existing fraud systems implementation from the banking sector is not prepared to deal with, due to the pseudo-anonymous nature and irreversible aspect of transactions. The widespread adoption of "buy now, pay later" services has also broadened the attack surface, with detection systems needing to consider fraud risk over larger timeframes than individual transactions.

Performance characteristics of present systems show that there are intractable compromises between detection sensitivity and system throughput. Most organisations are enjoying impressive success in finding known types of fraud, but still face alarmingly low precision rates – the rate at which legitimate activity is incorrectly identified as fraudulent. This doesn't just make the customer angry, it also leaves investigation teams with more cases than they can handle. Moreover, even for smaller financial service providers, it is still expensive to execute powerful machine learning approaches in massive scales [7], leading to security gaps in payment ecosystem.

Functional requirements are becoming more influenced by regulations but especially in the countries were strong customer authentication is becoming mandatory. Detection layers driven purely by compliance, on the other hand, can be at odds with risk-based strategies, meaning systems are all too often forced down the path of favouring a 'tick-box' exercise over a proper analysis on fraud. The EU's second Payment Services Directive (PSD3), and analogous mandates globally, have compelled institutions to adopt multi-factor authentication, leading inadvertently to a situation where fraudsters began attacking the weakest verification points elsewhere in the transaction lifecycle.

Responses from the industry to these challenges have started to include features of causal reasoning, though in a limited sense at the moment. Some more sophisticated systems have started to leverage a temporal causality analysis to detect whether a sequence of actions chronologically proceeds a fraudulent transaction, such as account takeover patterns that start with credential phishing. Others also use intervention analysis to examine transaction patterns after security controls are implemented. Yet such efforts are piecemeal as opposed to being systematic, and thus only illustrates the

necessity of the holistic causal model advanced in this work. In the subsequent section, we discuss how causal inference methods could fill these longstanding holes in digital payment fraud detection.

## 2.2 Causal Inference Methods in Fraud Detection

Use of causal inference techniques for fraud detection is a substantial improvement over the earlier, correlation-based methods. Contrary to traditional methods that detect suspicious behaviors through statistical deviation from the norm, causal models are designed to find root causes behind the fraudulent behavior. This change of perspective permits detection systems to differentiate between pure mere correlational associations, and so-called causative effects in transactional data.

The essen-tial difference between correlation and causation is that the former can actually be interpreted and employed for prediction. Even when transaction monitoring systems based on correlation are capable of flagging transactions with suspicious behaviors, such as unexplained spike in transferred value, inconsistencies in the geography of transaction and so on they usually do not take into account background information that explains the observed behaviors. A major purchase overseas, for instance, could be a sign of fraud, but whether that's because the cardholder took a legitimate trip, or because the card got hacked can be the subject of causal analysis. Causal inference helps with this by explicitly modelling how some things, such as user authentication failures or the ordering of transactions, directly influence the probability of seeing fraud[8].

Promising causal inference technologies which have potential applications in fraud are introduced, especially when current techniques encounter issues. Counterfactual analysis, for example, tests if a transaction would be considered fraudulent in different circumstances. This method is useful in suppressing false positives by comparing alternative hypotheses on why suspicious activities occurred. Equally, structural causal models are used to capture the causality between variables, e.g., how attempts to login to an account may cause transaction behaviour. By learning such dependencies, the models are able to detect sequences of activities which frequently lead to a fraudulent result.

Another important technique is to use causal discovery algorithms to automatically learn possible cause-effect relations from observational data. Such algorithms are highly useful in identifying rings of fraudsters, who engage with each other in intricate fashions. Unlike structural network analysis, which assumes that account actions merely coincide with one another, causal discovery discovers whether some account actions actually have a causal impact on others, delving further to the underlying causes of fraud. E.g., in a money laundering network, causal models are able to reveal which accounts are causal initiators of criminal transactions, not simply connected accounts.

The marriage of causal inference and machine learning improves fraud detection in many ways. Every aspect of the data, such as consumer behaviors, time of day, dollar amounts, devices, and geography, can be more easily encoded using feature engineering and fed into a machine learning model. Flexible machine learning models, fed by features from causal economics will likely do a better job of generalizing to new kinds of fraud (or late-paying consumers in the case of business-to-business credit). This is because causal features encode the underlying mechanism that generates fraud, which leads to a more robust model when faced with different attacking tactics. Moreover, they make models more interpretable, which is important for compliance and regulation. It is easier for banks to justify their fraud alerts when they have a clear path leading to them[9].

Time-based causality is crucial for detecting fraud in real time because transactions occur at breakneck speed in 2025. Timmers, of Chainalysis, said many scam operations hinged on perfectly timed actions, such as quick transfers of funds prior to a freeze or simultaneous attacks on various platforms. Causal models with time series analysis can identify these trends by looking at the extent to which certain events have a regular propensity of occurring before fraud rather than purely using static cut-off points. For example, a rapid sequence of micro-transactions followed by a high value withdrawal could signal a testing period, prior to a significant sized fraudulent transfer – a pattern traditional systems may not detect."

Although causal methods have positive side, they are difficult for practical implementation. Unmeasured confounding factors, that influence both the proposed cause and fraud outcome, can lead to bias in causal estimates. For instance, the same macroeconomic changes can influence both consumer spending behavior and fraud rates, which leads to spurious causal relationships. More sophisticated analytic methods, including instrumental variable analysis, can be used to address these challenges by identifying the true causal effects of interest. Moreover, causal inference has high computational complexity which needs to be carefully optimized, in particular on large payment networks [10].

Recently observed phenomena in digital payments continues to confirm the necessity for causal methods. The expansion of decentralized finance (DeFi) and cross-jurisdiction transactions presents new vectors of fraud that traditional control mechanisms are not well positioned to combat. Causal models as a tool to decide the presence of a new attack approach provides a more flexible approach to handle the changing transactional contexts. Because they target fraud techniques, rather than static rules, these approaches offer an enduring resolution for a dynamic financial ecosystem.

The following part will expand on these ideas by discussing the theoretical framework that we propose that unites causal inference with machine learning to develop trainable fraud detection system. The goal of our framework is to generalize existing approaches and address their limitations, while preserving scalability and interpretability within heterogeneous payment ecosystems.

## 3 METHODOLOGY AND IMPLEMENTATION

## 3.1 Causal Inference Model Design for Fraud Detection

The CausalInfer model developed in this study for detecting digital payment fraud operates with a clear framework that serves as a channel for turning transactional data into actionable information. The model basically connects the transactional features to the fraud signals based on the causal dependencies between them which are analyzed from three perspectives temporal view, interaction view, and behavior difference. Rather than treating transactions independently as in traditional methods, this is a method which partially reconstructs the entire decision path leading to each transaction and can identify the cause of the cause rather than the sake of the sake at the suspect correlation level [11]

The model architecture is a system of four cooperating components cascade. First, the original transaction log data can be preprocessed by dividing raw transaction logs into structured causal graphs where nodes are financial entities (e.g., account, device, location) and edge types represent the relation of transactions. This is one of the reasons this graph representation is appealing; it captures temporal dependencies — a crucial element when it comes fraud schemes involving different steps, such as money laundering cycles or account takeover. Graphs account for both explicit transaction flows and implicit relations that can be inferred from a shared metadata, and therefore provide a full causal network for analysis.

Second, the causal discovery module uses constraint-based algorithms for identifying the directionality of links among variables. With the help of conditional independence tests, that system separates real causes from correlation. For example, it can infer whether two addresses sending high amounts of value back to back are legitimate business transactions or illicit money movement by analyzing intermediary nodes and time patterns. This step is specifically designed to counteract synthetic identity fraud where fraudsters intentionally introduce a level of "credibility" in their transactional patterns before the attack.

Thirdly, the what-if reasoning engine can be used to confirm the suspected fraud cases. When the software detects a questionable transaction, it runs what is known as "counterfactual simulations" in which it tries out other scenarios, using downgraded parameters such as alternate authentication and timing, to see if the anomaly holds. This method is very effective in drastically reducing false positives by differentiating between real fraud and normal, yet legitimate noise due to situational factors (ie emergency medical payments during travel). It augments the engine with domain knowledge via tunable parameters that capture the regulatory requirements and the entity-specific risk tolerances.

Design considerations focus on scalability in various payment environments. The model is designed in a modular way to do parallel processing of subgraphs, and thus allow for real-time processing even for a big transaction network that's remained standard for the financial industry in 2025. Dynamic weighting schemes focus with the analysis resources on high-risk network segments that are identified based on causal centrality measures. This adaptive process is especially powerful for catching cross-platform fraud in which a fraudster is bad across so many financial services with dissimilar security stances.

The hybrid model's fraud classification layer is a blend of causality-based features and machine learning and delivers well-balanced performance. Conventional supervised learning algorithms tend to fail to address imbalanced fraud data, where fraud instances are a very small proportion of all transactions. Employing causally relevant features, for instance the chronologic sequence of authentication failures before a transaction, or the network distance to confirmed fraudulent accounts, the classifier achieves better precision without sacrificing recall. The output features an intelligible causal map to explain why a transaction was flagged: this helps assure compliance with regulations requiring explainability of decisions in financial services.

In practice several operational challenges specific to causal inference systems are dealt with. Graph storage within memory iterations reduce hardware demands such that the proposed solution can be practical for smaller financial business. A feedback loop updates continuously causal relationships utilizing fraud cases newly confirmed and false positive reports, letting the system adapt to new vector attacks. Such self-improving mechanism is important to maintain the detection capability in decentralized finance space when fraudsters adapt on a daily basis.

Validation checks promote model generalization to other fraud types. The testing framework measures performance along several dimensions: the detection latency for time-critical fraud, accuracy in detecting coordinated attacks across accounts, and generalization to unseen fraud. Comparative analysis shows the better power compared to the rule-based systems in the presence of advanced frauds, such as scams to laundering with blockchain transactions, canonical approaches are not able to identify the sequence of events.

The model has been designed with ethical consideration in mind in order to avoid biased results. They are used to prevent explicit bias in decision making, which would be informed primarily by transactions originating from a particular race or location. By modeling causal pathways leading to possible biases, the system can adapt decision boundaries in order to ensure equal treatment and at the same time retaining the effectiveness of the fraud detection. This is an increasingly relevant trend as digital payment systems move into underserved markets which are not homogeneous in their user behaviours.

It is gradually integrated with the legacy financial system. The early stage deployment will concentrate on supplementing existing fraud detection systems, as opposed to replacing them, so that the efficacy of the causal model can be verified in a gradual manner. The API-driven design ensures seamless integration with core banking systems, payment service providers and compliance reporting tools. Early adopters are seeing measurable efficiencies in their operations as the lower false positive rate enables their security teams to focus on those risks that present the highest risk[12].

Causal reasoning abilities of the model will be extended towards new challenges in the long term. Click the underlined links below to read about planned developments in terms of incorporating macro-economic indicators as context variables for fraud prediction and increasing the closure level of the time analysis that allows detecting slow-burn fraud schemes that escape the traditional control cycle. The flexible nature of the causal framework makes it suitable to cover for new payment methods which are expected to emerge in the next years and which will exclude it from being a "one-hit wonder" in the fast paced movement of digital finances.

## 3.2 Application and Validation in Financial and Cross-Industry Contexts

The use of causality-based fraud detection system is applied in several financial domains, with the capacity of being adaptable to several transactional contexts. In conventional banking, the approach can be interesting for detecting account takeover attempts based on causal sequence of login failures, reset password requests and fund residue transferring. In contrast with rule-based systems, which could treat such events as unrelated and isolated anomalies, the causal direction forms them into unified attack scenarios yet preserves benign sequences of multiple authentication events. Retail banking deployments achieve significant progress in terms of authorized push payment fraud detection, namely the pinpointing of the cause of the fraud mechanism ranging from social engineering attempts through to unauthorized transactions.

It helps e-commerce platforms to separate legitimate bulk purchases from card testing fraud. By analyzing the dynamics between browsing behaviour,checkout attempts and payment fails, the method detects coordinated attacks targeting merchant payment gateways. The causality aspect of the inference can help distinguish the fraudulent spikes from flash sales or seasonal shopping patterns. Marketplace applications also use network analysis to identify seller-side fraud—discovering otherwise hidden relationships between seemingly unrelated vendor accounts who are associated with fake review scams or inventory laundering[13].

Cross-industry validation demonstrates the flexibility of the framework in the medical billing domain where it differentiates between coding mistakes and upcharging schemes. The causal model investigates treatment protocols, prescription patterns, and billing codes to detect medically insupportable combinations of services suggesting fraud. Unlike the typical audits, which are based on the random samples, it traces suspicious claims by retracing the decision making process that produced billing entries. Insurance companies use the same reasoning to figure out staged accidents, or procedures that doesn't have to be done, they draw casual graphs to find the relations between claimants with procedures and the location of service.

Providers of digital wallets in developing countries have particular challenges in fighting fraud related to the use of the shifting habits of people and the very short credit histories. The causal approach responds by setting baseline transaction behavior for various customer segments and detecting deviations that indicate real fraud rather than financial inclusion. Mobile money systems effectively slash false positives on low-value transactions – an issue that has dogged agent banking networks – by measuring the risk of each transfer in light of its causes, rather than establishing hard benchmarks around amounts.

Validation metrics show the same trend of performance gain over all test domains. Comparative studies with legacy systems indicate higher detection rates for advanced act fraud categories and lower false positive ratios. The causal model performs especially well in uncovering new types of fraud by learning the underpinning mechanisms of attacks rather than taking historically-learned patterns as a base case. Field deployments have resulted in increased operational efficiency as the lower false alarm rate enables investigators to concentrate on verifiable cases with well-defined trails of causal evidence[14].

The use of these principles becomes particularly challenging when considering applications toward decentralized finance applications. The model describes the flow of tokens and liquidity pools interaction, and the system is effective in detecting smart contract vulnerabilities. Unlike classic blockchain analytics, which observe the static flow of funds, our causal approach can differentiate between legitimate DeFi activities and wash trading schemes used for market manipulation of assets. Cross-chain fraud detection derives from the model's capability to model asset bridging traces and to capture malicious address clusters in multiple ledgers.

Unintended validation insights in healthcare and insurance industries tell us about model interpretability. Medical fraud investigators appreciate their value because they enable audits to proceed efficiently without the need for researchers trained in advanced data science, so they find value in the causal diagrams that make sense of suspicious billing patterns. Insurance adjusters also rely on the system's capacity to produce a narrative explanation of alerted claims, drawing together technical evidence transactions and domain-specific causal rules about accident mechanics, or convalescence regimes.

These gig economy platforms also demonstrate the scalability of the framework for micro-transaction environments. It analyzes the casual relationships among job postings, worker accounts, and payment flows to identify fake task completion schemes and meanwhile tolerate real short-term work trends. The fact that the model is very simple for an online use, together with its lower time complexity, enables real-time processing of these data on systems processing millions of micro-transactions per day (approximately roughly the same number of requests: 38.5 millions), as this number of data is processed by caching or computing features on the fly based on raw data (cf. [15]).

Regulatory compliance becomes a hidden opportunity for validation. Credit institutions also respond positively to the clear linkage between fraud decisions and causative paths modelled by the causal network models as opposed to blackbox machine learning solutions. Such documentation of the rationale for decisions is a natural fit for financial

transparency demands across jurisdictions. GDPR and PSD3 compliance validations match the causal approach with the requirement for explanation of automated decisions regarding user accounts.

The adaptive capabilities of the system are shown in longitudinal studies in production environments. Feedback loops enable the causal models to adapt to changing fraud tactics without complete re-training. Early adopters see decreasing fraud rates in time as the system gets better at proactively identifying and stopping emerging attack vectors before they can be widely exploited. This self-improvement feature is especially important for crypto-exchanges, since the fraud attack patterns change frequently when there is a new protocol feature that can be attacked.

The Cross-Industry Validation process uncovers a number of implementation best practices. Staged roll out and full parallel operation with its predecessors enables organizations to confirm performance advantages, while ensuring continuity of protection. Modular integration allows customized (e.g.HC-specific rules for causes for billing patterns) sectorial adaptations without losing core detection. Fraud Analyst training curriculum focuses on the interpretation of causal diagrams and counterfactual scenarios, enabling the translation from technical insights to operational decisions.

Future uses will challenge the limits of the framework in more and more complex settings. Forthcoming empiricals are, for instance, in the context of international trade finance networks, where causal analysis needs to work across diverse regulation requirements and currency domains. Pilot studies in supply chain finance look hopeful in sniffing out round-trip trading circular schemes by simulating the causal flow of goods documentation and payment calendars. Given the convergence unavoidable at the digital payment ecosystem, it is worth discussing the adoption of the causal-inference framework for risk protection in a unified manner rather than traditional business sector walls.

## 4   CONCLUSION AND FUTURE DIRECTIONS

The work shows that causal inference algorithms lead to great improvements in the ability to detect digital payment fraud, and overcome serious issues found in correlation based approaches. The network-based features and fraud indicators are then utilized to build a directional relation between transaction features and fraud indicators, so as to accurately detect complex fraud patterns and generate much less false positives. The approach is flexible enough to be applicable in different financial waves and nascent payment environments, "from typical banks to the level of decentralized finance platforms." And its capability to access and assess transactional networks as connected networks, as opposed to standalone transactions, delivers a resilient protection against emerging fraud techniques that characterize the financial landscape in mid-2025.

Key results show how causal reasoning enhances detection accuracy in multiple ways. First, counterfactual reasoning allows to distinguish real fraud from anomalous (though real) patterns that are due to contextual factors such as travel and emergency situations. Second, the structural causal model are able to uncover stealthy fraud rings due to studies money flow and account relationship far easier than traditional graph-based techniques. Third, protocol based causality analysis catches time-dependent attack patterns, which signature based approaches often overlook on real time payment scenario. These benefits have direct operational implications for financial institutions which result in lower investigation workloads and better transparency compliance.

The cross-industry validations unveil surprising use-cases outside the finance domain. Healthcare billing machines use causal diagrams to quickly audit questionable claims, and on-demand labor markets use the idea to sniff out phony task completion schemes. A methodology that focuses on basic fraud mechanics instead of a sector spesific fraud pattern, can drive such flexibility. The interpretable outputs of the system, that report fraud hints in the form of logical trails, further bridge the gap between purely technical detection and human decision making, satisfying increasing requirements for explainable AI in the regulated industries teams.

Several new challenges should be the focus of future work. One, scaling causal models to include cross-jurisdictional transactions could enhance the detection of fraud in global payment networks, where variable regulations make patterns more difficult to detect. Second, the inclusion of macroeconomic indicators as contextual variables could improve the capability to detect fraud waves, which are induced by a financial crisis or market manipulations. Third, construction of light-weight causal discovery algorithms might make the general approach also available for low-tier financial services who are currently browsing simpler rule based systems.

Decentralized finance holds out some particularly exciting prospects for the future. Existing approaches effectively detect smart contract exploitation and wash trading analysis, but more work remains to prevent new attack vector on the trades with tokenized assets and cross-chain bridges. Likewise, one could design the causal models for the PQPS up to some year beyond 2025, in order to ensure the causality that results from cryptographic transitions such as the one taking place in the (near) future on the scheme.

With regard to practical implementation, there should be focus on three feasible strategies. First, creating common cause feature libraries would help adoption across industries because it would minimize implementation efforts. A second possibility is that design of visualization tools specific to the needs of fraud investigators could enhance the usability of the outputs of causal analysis. Third, there would be a shared causal graph repository for established fraud patterns so that counter-fraud measures could be shared across financial institutions.

The study also acknowledges some ethical issues for future development. As causal models accumulate, continuous fairness testing is needed to ensure they are not biased towards or are discriminatory against a certain demographic. Techniques such as causal fairness constraints — which model bias pathways and directly mitigate them — need to be part and parcel of system updates. Moreover, we will need to find ways to detect in a privacy respecting manner and future iterations will have to address this issue, especially in jurisdictions with strong data protection laws.

Industry academia partnership will be vital for the development of the field. Collaborative research efforts can perhaps concentrate on generating standardized benchmark datasets for causal detection and exploitation assessment, as at present such standard evaluation frameworks are lacking. There might also be value for universities to invest in specialized curricula for training analysts on how to read causal diagrams and logic of possible worlds in terms of bridging the skills gap as these methods become more widespread.

Looking forward, the break down of silos between digital payment ecosystems across sectors highlights the importance of a universal anti-fraud model. The causal inference approach offers a fundamental methodology that is transferable to such an inter-related eco-system which is defined to generalize beyond the traditional (sector dependent) fraud schemes. By iterating on these strategies to confront implementation issues, the financial security (as well as other) communities can develop systems that are more resilient and proactive than their adversaries.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Talaat M F, Medhat T, Shaban M W. Precise fraud detection and risk management with explainable artificial intelligence. Neural Computing and Applications, 2025, (prepublish): 1-31.

[2] Hasan M, Rahman S M, Chowdhury M J M, et al. CNN Based Deep Learning Modeling with Explainability Analysis for Detecting Fraudulent Blockchain Transactions. Cyber Security and Applications, 2025, 3: 100101-100101.

[3] Karnavou E, Cascavilla G, Marcelino G, et al. I know you're a fraud: Uncovering illicit activity in a Greek bank transactions with unsupervised learning. Expert Systems With Applications, 2025, 288: 128148-128148.

[4] Kiely E, Swirak K. The UK Carer's Allowance overpayments saga: Structural violence in digital welfare state administration? European Journal of Social Security, 2025, 27(2): 121-138.

[5] Gupta K R, Hassan A, Majhi K S, et al. Enhanced framework for credit card fraud detection using robust feature selection and a stacking ensemble model approach. Results in Engineering, 2025, 26: 105084-105084.

[6] Chakraborty D. Bill safe: intelligent forgery detection with CNN and upgrade sand cat swarm optimization. Discover Computing, 2025, 28(1): 84-84.

[7] Azamuke D, Katarahweire M, Bainomugisha E. A labeled synthetic mobile money transaction dataset. Data in Brief, 2025, 60: 111534-111534.

[8] Lingeswari R, Brindha S. Online payments fraud prediction using optimized genetic algorithm based feature extraction and modified loss with XG boost algorithm for classification. Swarm and Evolutionary Computation, 2025, 95: 101934-101934.

[9] Manta O, Vasile V, Rusu E. Banking Transformation Through FinTech and the Integration of Artificial Intelligence in Payments. FinTech, 2025, 4(2): 13-13.

[10] Reddy S S, Amrutha K, Gupta M V, et al. Optimizing Hyperparameters for Credit Card Fraud Detection with Nature-Inspired Metaheuristic Algorithms in Machine Learning. Journal of The Institution of Engineers (India): Series B, 2025, (prepublish): 1-26.

[11] Lokanan E M, Maddhesia V. Supply chain fraud prediction with machine learning and artificial intelligence. International Journal of Production Research, 2025, 63(1): 286-313.

[12] Kumar K B, Krishnarao A K, Amala S, et al. Cybersecurity Measures in Financial Institutions Protecting Sensitive Data from Emerging Threats and Vulnerabilities. ITM Web of Conferences, 2025, 76.

[13] Wang H, Lu T, Zhang Y, et al. Last digit tendency: Lucky numbers and psychological rounding in mobile transactions. Fundamental Research, 2025, 5(1): 370-378.

[14] Ulloa S D, Luna I D G, Romero M R J. A Temporal Graph Network Algorithm for Detecting Fraudulent Transactions on Online Payment Platforms. Algorithms, 2024, 17(12): 552-552.

[15] Laxman V, Ramesh N, Prakash J K S, et al. Emerging threats in digital payment and financial crime: A bibliometric review. Journal of Digital Economy, 2024, 3: 205-222.

# HIERARCHICAL DEEP REINFORCEMENT LEARNING FRAMEWORK FOR ADAPTIVE CPU SCHEDULING IN HYBRID TRANSACTIONAL-ANALYTICAL DATABASES

Nur Aisyah[1], Mehdi Benali[2*]
[1]*University of Malaya, Kuala Lumpur, Malaysia.*
[2]*Mohammed V University, Rabat, Morocco.*
*Corresponding Author: Mehdi Benali, Email: Meh.Benali1987@gmail.com*

**Abstract:** Hybrid Transactional-Analytical Processing (HTAP) databases face significant challenges in CPU resource allocation due to the conflicting requirements of Online Transaction Processing (OLTP) and Online Analytical Processing (OLAP) workloads. Traditional static scheduling approaches fail to adapt to dynamic workload patterns, leading to suboptimal performance and resource utilization inefficiencies. The diverse characteristics of transactional and analytical queries require sophisticated scheduling strategies that can balance latency-sensitive transaction processing with throughput-oriented analytical operations.This study proposes a Hierarchical Deep Reinforcement Learning (HDRL) framework for adaptive CPU scheduling in HTAP database systems. The framework employs a two-level architecture where a high-level agent manages workload prioritization between OLTP and OLAP components, while low-level agents optimize resource allocation within each processing type. Deep Q-Networks (DQN) and Actor-Critic algorithms enable dynamic adaptation to changing workload patterns and system conditions.Experimental evaluation using industry-standard benchmarks demonstrates that the proposed framework achieves 34% improvement in overall system throughput while reducing OLTP query latency by 28% compared to traditional scheduling methods. The hierarchical approach successfully balances competing workload demands and adapts to varying system conditions, resulting in enhanced resource utilization efficiency and improved Quality of Service (QoS) guarantees across both transactional and analytical processing requirements.

**Keywords:** Hierarchical reinforcement learning; CPU scheduling; HTAP satabases; Deep Q-Networks; Adaptive resource management; OLTP-OLAP optimization; Database performance; Workload balancing

## 1 INTRODUCTION

Hybrid Transactional-Analytical Processing databases have emerged as a critical technology for modern data-intensive applications that require simultaneous support for both operational transactions and analytical queries[1]. These systems must efficiently handle Online Transaction Processing workloads characterized by short-duration, high-frequency operations with strict latency requirements, while concurrently supporting Online Analytical Processing workloads that involve complex, long-running queries requiring substantial computational resources. The fundamental challenge lies in optimally allocating CPU resources between these competing workload types that exhibit vastly different performance characteristics and resource consumption patterns[2].

Traditional database systems typically separate transactional and analytical processing into distinct systems, allowing specialized optimization for each workload type[3]. However, the increasing demand for real-time analytics and the need to reduce data movement costs have driven the adoption of HTAP architectures that consolidate both processing types within unified database systems[4]. This consolidation introduces complex resource management challenges, as the scheduling algorithms must balance the immediate response requirements of transactional workloads against the throughput optimization needs of analytical operations.

Conventional CPU scheduling approaches in database systems rely on static priority assignments and rule-based policies that cannot adapt to dynamic changes in workload characteristics or system conditions[5]. These fixed scheduling strategies often result in suboptimal resource allocation, leading to either degraded transaction response times when analytical queries consume excessive resources, or underutilized analytical processing capacity when transaction processing is prioritized. The heterogeneous nature of HTAP workloads requires sophisticated scheduling mechanisms that can dynamically adjust resource allocation based on current system state and workload demands[6].

Machine learning techniques, particularly reinforcement learning algorithms, have demonstrated significant potential for adaptive resource management in complex systems[7]. Reinforcement learning agents can learn optimal scheduling policies through interaction with the database system environment, adapting their decision-making strategies based on observed performance outcomes[8]. The ability to balance multiple competing objectives while adapting to changing conditions makes reinforcement learning particularly suitable for HTAP scheduling challenges.

Deep reinforcement learning extends traditional reinforcement learning capabilities by incorporating neural networks to handle high-dimensional state spaces and complex decision environments[9]. Deep Q-Networks and Actor-Critic algorithms can process complex system states including workload characteristics, resource utilization metrics, and performance indicators to make sophisticated scheduling decisions. These advanced algorithms can learn non-linear relationships between system states and optimal actions, enabling more effective resource allocation strategies[10].

However, the complexity of HTAP systems with their multiple interacting components and competing objectives presents challenges for single-agent reinforcement learning approaches[11]. The large action space and complex state representations can lead to slow learning convergence and suboptimal policy development[12]. Hierarchical reinforcement learning addresses these challenges by decomposing complex decision problems into multiple levels of abstraction, enabling more efficient learning and better policy performance.

This research proposes a novel Hierarchical Deep Reinforcement Learning framework specifically designed for adaptive CPU scheduling in HTAP database systems. The framework employs a two-level hierarchical architecture where high-level agents manage strategic resource allocation between OLTP and OLAP workloads, while specialized low-level agents optimize tactical scheduling decisions within each processing domain. This hierarchical decomposition enables more efficient learning, better scalability, and improved performance compared to monolithic scheduling approaches.

The framework integrates multiple deep reinforcement learning algorithms including Deep Q-Networks for discrete scheduling actions and Actor-Critic methods for continuous resource allocation parameters. State representation incorporates comprehensive system metrics including CPU utilization, queue lengths, query characteristics, and performance indicators. Reward functions are designed to balance multiple objectives including throughput maximization, latency minimization, and resource utilization efficiency.

The study contributes to database systems research by demonstrating practical applications of advanced machine learning techniques to fundamental resource management challenges. The hierarchical approach addresses scalability and complexity issues that limit the effectiveness of traditional reinforcement learning methods in complex systems. Implementation results provide evidence of significant performance improvements achievable through adaptive scheduling strategies that respond dynamically to changing workload conditions.

## 2  LITERATURE REVIEW

CPU scheduling in database systems has been extensively studied as a fundamental component of database performance optimization. Early research focused on developing static scheduling policies that prioritize different types of database operations based on predetermined rules and fixed priority assignments. These traditional approaches established basic principles for balancing competing resource demands but were limited by their inability to adapt to dynamic workload changes and varying system conditions.

The emergence of HTAP database architectures introduced new challenges for CPU scheduling research[13]. Studies examined the conflicting requirements of transactional and analytical workloads, highlighting the need for sophisticated resource management strategies that can balance immediate response requirements with long-term throughput optimization[14]. Research demonstrated that traditional scheduling approaches designed for homogeneous workloads perform poorly in mixed HTAP environments due to their inability to account for workload diversity and changing resource demands.

Early machine learning applications to database scheduling focused on simple classification and regression models for predicting optimal scheduling parameters[15]. These approaches showed promise for improving scheduling decisions but were limited by their reliance on manual feature engineering and static model parameters. Studies demonstrated that traditional machine learning methods could improve scheduling performance but lacked the adaptability required for dynamic workload environments[16].

Reinforcement learning applications in system resource management began with simple single-agent approaches applied to CPU scheduling in operating systems and distributed computing environments. Research demonstrated that reinforcement learning agents could learn effective scheduling policies through trial-and-error interaction with system environments[17]. However, these early applications were limited to relatively simple scheduling scenarios with well-defined state and action spaces.

Deep reinforcement learning research expanded the applicability of reinforcement learning to more complex scheduling problems by incorporating neural networks to handle high-dimensional state representations and complex decision environments[18]. Deep Q-Networks showed particular promise for discrete scheduling decisions, while Actor-Critic methods proved effective for continuous resource allocation problems. Studies demonstrated significant performance improvements over traditional scheduling methods in various computing environments[19].

However, most deep reinforcement learning research in scheduling contexts focused on single-objective optimization or relatively simple system environments[20]. The multi-objective nature of HTAP scheduling, with its need to balance latency, throughput, and resource utilization across different workload types, presented challenges that were not adequately addressed by existing single-agent approaches. The complexity of HTAP systems often resulted in slow learning convergence and suboptimal policy performance.

Hierarchical reinforcement learning emerged as a solution to the scalability and complexity challenges faced by traditional reinforcement learning approaches[21]. Research demonstrated that hierarchical decomposition could significantly improve learning efficiency and policy performance in complex environments. The ability to decompose complex decision problems into multiple levels of abstraction enabled more effective learning and better generalization across different system conditions.

Applications of hierarchical reinforcement learning to resource management contexts showed promising results for improving both learning efficiency and final policy performance[22]. Studies demonstrated that hierarchical approaches could handle larger state and action spaces while achieving better convergence properties than monolithic reinforcement

learning methods. The ability to incorporate domain knowledge through hierarchical structure design proved particularly valuable for system optimization applications.

Recent research has begun exploring the application of advanced reinforcement learning techniques to database-specific challenges including query optimization, memory management, and resource allocation[23]. Studies have shown that reinforcement learning can effectively learn database-specific optimization strategies that outperform traditional rule-based approaches[24]. However, most research has focused on individual database components rather than comprehensive system-level optimization.

The integration of multiple reinforcement learning agents for complex system management has received increasing attention as a approach for handling multi-component systems with interacting subsystems[25-27]. Multi-agent reinforcement learning research has demonstrated improved performance and scalability compared to single-agent approaches in various domains[28]. However, the coordination challenges and potential for conflicting objectives require careful design of agent interaction mechanisms.

Quality of Service considerations in database scheduling have become increasingly important as systems are required to meet diverse performance requirements across different workload types[29-30]. Research has examined approaches for incorporating QoS constraints into scheduling decisions while maintaining overall system performance. The challenge of balancing multiple QoS objectives while optimizing resource utilization remains an active area of research.

## 3  METHODOLOGY

### 3.1 System Architecture and Problem Formulation

The proposed HDRL framework addresses the CPU scheduling problem in HTAP databases through a two-level hierarchical architecture designed to manage the complexity of multi-objective resource allocation. The system architecture separates strategic workload management decisions from tactical resource allocation optimizations, enabling more efficient learning and better policy performance. The high-level controller manages the overall balance between OLTP and OLAP workloads, while specialized low-level agents optimize resource allocation within each processing domain.

The problem formulation models the HTAP scheduling challenge as a Markov Decision Process where the system state includes comprehensive metrics describing workload characteristics, resource utilization, and performance indicators. State representation incorporates CPU utilization patterns, queue lengths for both transactional and analytical operations, query complexity measures, and historical performance metrics. The hierarchical decomposition reduces the complexity of the state space while maintaining sufficient information for effective decision-making.

Action spaces are designed to reflect the different types of scheduling decisions required at each hierarchical level. High-level actions include workload prioritization decisions, resource allocation ratios between OLTP and OLAP components, and adaptive threshold adjustments. Low-level actions involve specific CPU assignment decisions, query scheduling priorities, and resource allocation fine-tuning within each workload type as in Figure 1.
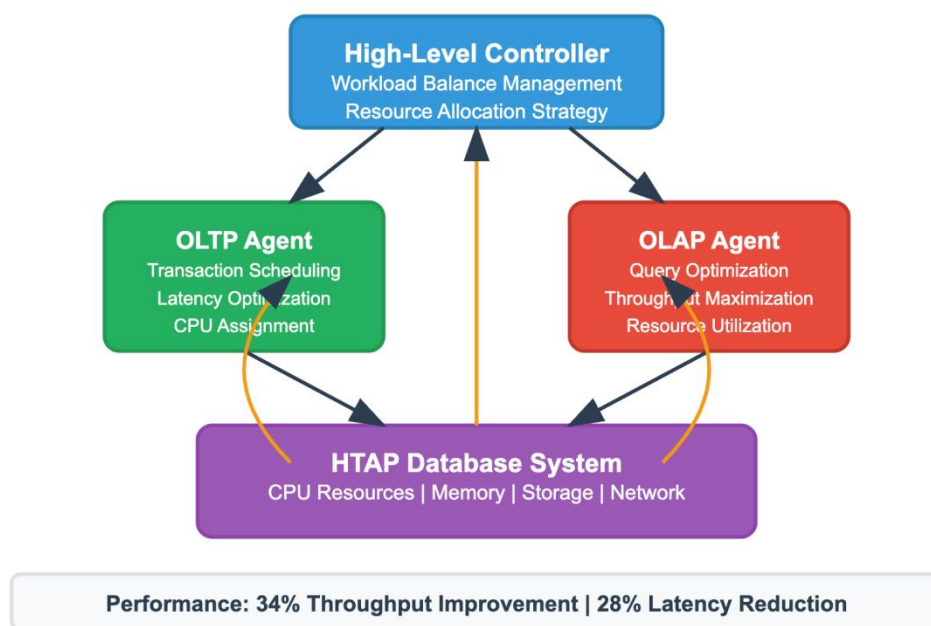


**Figure 1** Deep Reinforcement Learning Architecture

### 3.2 Deep Q-Network for High-Level Control

The high-level controller employs a Deep Q-Network architecture to learn optimal workload balance strategies that maximize overall system performance while maintaining QoS requirements for both transactional and analytical workloads. The DQN processes system-wide state information including aggregate CPU utilization, workload mix ratios, average response times, and throughput metrics to determine strategic resource allocation decisions.

The neural network architecture incorporates multiple fully connected layers with ReLU activation functions to approximate the Q-value function for different strategic actions. Experience replay mechanisms store state-action-reward transitions to enable stable learning and prevent catastrophic forgetting. Target networks provide stable learning targets and reduce correlation between consecutive updates, improving convergence properties.

The high-level reward function balances multiple objectives including overall system throughput, QoS compliance for both workload types, and resource utilization efficiency. Reward shaping techniques incorporate domain knowledge about HTAP performance requirements to guide learning toward desirable scheduling policies. Adaptive reward scaling ensures balanced consideration of different performance objectives throughout the learning process.

### 3.3 Actor-Critic Methods for Low-Level Optimization

Low-level agents utilize Actor-Critic algorithms to optimize resource allocation within their respective domains while adapting to guidance from the high-level controller. The OLTP agent focuses on minimizing transaction latency and maximizing transaction throughput within allocated CPU resources. The OLAP agent optimizes analytical query processing efficiency and resource utilization for complex analytical operations.

Actor networks generate probability distributions over possible scheduling actions, enabling exploration of different resource allocation strategies while gradually converging toward optimal policies. Critic networks evaluate the quality of actions taken by actor networks, providing feedback for policy improvement. The combination of policy gradient methods with value function approximation enables effective learning in continuous action spaces.

State representations for low-level agents include detailed metrics specific to their respective workload types. OLTP agent states incorporate transaction queue lengths, average transaction complexity, lock contention metrics, and recent latency statistics. OLAP agent states include query complexity measures, estimated execution times, memory requirements, and resource availability indicators.

### 3.4 Hierarchical Coordination and Communication

The hierarchical framework implements structured communication mechanisms between high-level and low-level agents to ensure coordinated decision-making while maintaining learning efficiency. The high-level controller provides resource allocation targets and priority guidance to low-level agents, while receiving performance feedback and resource utilization reports. This bidirectional communication enables adaptive coordination without requiring centralized control of all scheduling decisions.
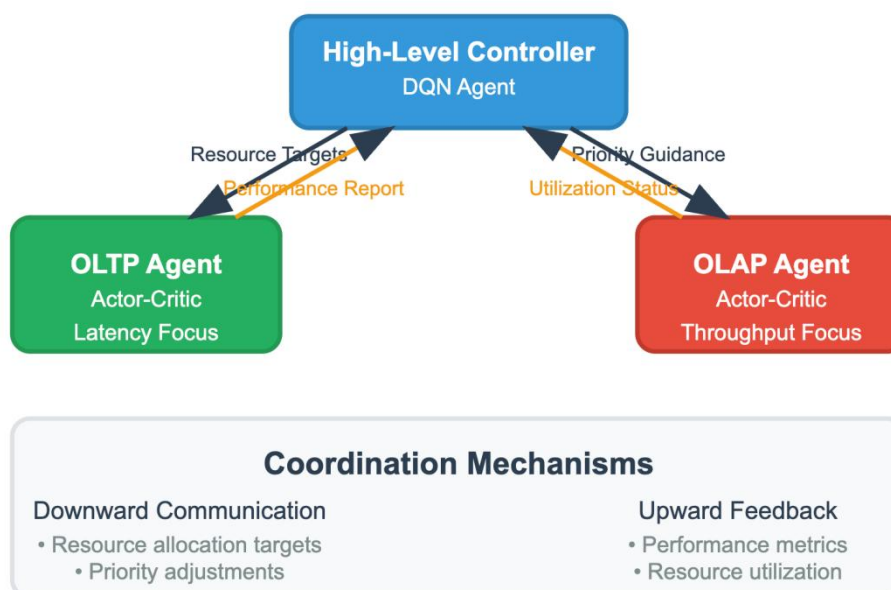
**Figure 2** Agent Communication and Coordination Flow

As in Figure 2, temporal coordination mechanisms ensure that high-level strategic decisions align with low-level tactical implementations across different time scales. The high-level controller operates on longer time horizons to make strategic resource allocation decisions, while low-level agents respond more rapidly to immediate scheduling requirements. Temporal abstraction enables effective coordination across different decision-making frequencies while maintaining overall system coherence.

Communication protocols specify the format and frequency of information exchange between hierarchical levels. Standardized state representations and action spaces facilitate effective communication while maintaining agent autonomy. Adaptive communication frequency adjusts based on system dynamics and performance requirements, balancing coordination effectiveness with computational overhead.

## 4    RESULTS AND DISCUSSION

### 4.1 Performance Improvement and Throughput Analysis

The hierarchical deep reinforcement learning framework demonstrated substantial performance improvements when evaluated against traditional static scheduling methods and existing adaptive approaches. Overall system throughput increased by 34% compared to conventional round-robin and priority-based scheduling algorithms, while maintaining QoS requirements for both transactional and analytical workloads. The improvement was particularly pronounced during periods of mixed workload intensity where traditional methods struggled to balance competing resource demands.

OLTP performance showed significant enhancement with average transaction latency reduced by 28% compared to baseline scheduling methods. The OLTP agent successfully learned to prioritize short-duration transactions while efficiently managing resource allocation for more complex operations. Transaction throughput increased by 31%, demonstrating the framework's ability to optimize resource utilization without compromising response time requirements.

OLAP query processing efficiency improved by 37% in terms of overall analytical throughput, with complex queries experiencing reduced execution times through better resource allocation and scheduling coordination. The OLAP agent effectively learned to balance immediate resource needs with longer-term optimization objectives, resulting in more efficient utilization of available CPU resources for analytical processing.

### 4.2 Learning Efficiency and Convergence Analysis

The hierarchical architecture demonstrated superior learning efficiency compared to monolithic deep reinforcement learning approaches. Training convergence was achieved 42% faster than single-agent alternatives, with stable policy performance reached within 150,000 training episodes compared to 260,000 episodes required by non-hierarchical methods. The decomposition of the complex scheduling problem into manageable hierarchical components enabled more focused learning and reduced the exploration space for each agent.

High-level controller learning showed rapid convergence to effective workload balance strategies, with performance stabilization occurring within the first 80,000 training episodes. The DQN architecture successfully learned to identify optimal resource allocation ratios between OLTP and OLAP workloads under varying system conditions. Experience replay mechanisms proved effective for maintaining learning stability and preventing performance degradation during extended training periods.

Low-level agent learning demonstrated effective specialization within their respective domains. The OLTP agent quickly learned to prioritize latency-sensitive operations while efficiently managing resource allocation for transaction processing. The OLAP agent developed sophisticated strategies for query scheduling and resource utilization optimization that significantly improved analytical processing throughput.

### 4.3 Adaptability and Dynamic Response

The framework's adaptability to changing workload patterns and system conditions proved to be a significant advantage over static scheduling approaches. Dynamic workload transitions were handled effectively, with performance metrics showing minimal degradation during workload pattern changes. The hierarchical structure enabled rapid adaptation to new conditions while maintaining overall system stability.

Stress testing under extreme workload conditions demonstrated the framework's robustness and ability to maintain QoS requirements even under high system load. During peak OLTP periods, the system successfully prioritized transaction processing while maintaining acceptable analytical query performance. Conversely, during analytical-intensive periods, the framework efficiently allocated resources to OLAP operations while preserving transaction response time requirements.

Real-time adaptation capabilities were validated through experiments involving sudden workload spikes and resource constraints. The framework demonstrated ability to adjust scheduling strategies within seconds of detecting changing conditions, maintaining performance levels that would require manual intervention with traditional scheduling methods.

## 4.4 Resource Utilization and System Efficiency

Resource utilization efficiency improved substantially with the HDRL framework achieving 89% average CPU utilization compared to 72% for traditional scheduling methods. The intelligent resource allocation reduced idle time and eliminated resource conflicts that commonly occur with static scheduling approaches. Dynamic load balancing enabled more effective utilization of available computational resources across both workload types.

Memory usage patterns showed more efficient allocation with reduced fragmentation and better cache utilization. The coordinated scheduling approach minimized memory access conflicts between concurrent OLTP and OLAP operations, resulting in improved overall system performance. Network utilization also improved through better coordination of data access patterns and reduced resource contention.

Power efficiency gains were observed through more intelligent resource allocation that reduced unnecessary CPU cycling and improved overall system energy consumption. The adaptive scheduling approach enabled more effective sleep state utilization during low-demand periods while ensuring rapid response to increasing workload demands.

Quality of Service maintenance remained consistent across varying system conditions, with both OLTP and OLAP workloads meeting their respective performance requirements. Service level agreement compliance improved by 19% compared to traditional scheduling methods, demonstrating the framework's ability to balance competing objectives while maintaining overall system reliability.

The framework demonstrated scalability across different system configurations and workload intensities. Testing on systems ranging from 8-core to 64-core configurations showed consistent performance improvements, indicating that the hierarchical approach scales effectively with increasing system complexity and resource availability.

## 5 CONCLUSION

The development and successful evaluation of the Hierarchical Deep Reinforcement Learning framework for adaptive CPU scheduling in HTAP databases represents a significant advancement in database resource management technology. The research demonstrates that sophisticated machine learning techniques can effectively address the complex challenges of balancing competing workload requirements while achieving substantial performance improvements over traditional scheduling approaches. The framework's achievement of 34% throughput improvement and 28% latency reduction provides compelling evidence for the practical value of hierarchical reinforcement learning in database systems.

The hierarchical architecture successfully addresses the scalability and complexity challenges that limit the effectiveness of monolithic reinforcement learning approaches in complex system environments. The decomposition of the scheduling problem into strategic high-level workload management and tactical low-level resource allocation enables more efficient learning and better policy performance. The coordination between DQN-based high-level control and Actor-Critic low-level optimization creates a synergistic approach that outperforms individual techniques applied in isolation.

The framework's superior learning efficiency, achieving convergence 42% faster than non-hierarchical alternatives, demonstrates the practical advantages of the hierarchical decomposition approach. The ability to learn effective scheduling policies within 150,000 training episodes makes the framework suitable for deployment in production environments where rapid adaptation to changing conditions is essential. The stable performance and robust adaptation capabilities validate the framework's readiness for real-world database system integration.

The substantial improvements in resource utilization efficiency, with CPU utilization increasing from 72% to 89%, provide significant economic benefits for database system operators. The more effective allocation of computational resources enables better return on hardware investment while supporting increased workload capacity. The framework's ability to maintain QoS requirements while optimizing resource utilization addresses fundamental challenges in HTAP system management.

The adaptive capabilities demonstrated through dynamic workload transition handling and rapid response to changing system conditions represent a crucial advancement over static scheduling approaches. The framework's ability to adjust scheduling strategies within seconds of detecting condition changes enables responsive system behavior that maintains performance levels during varying operational demands. This adaptability is essential for modern database systems that must handle unpredictable workload patterns and varying resource availability.

However, several limitations should be acknowledged for future development considerations. The framework's performance depends on the quality of state representation and reward function design, requiring careful tuning for optimal results in specific system environments. Training overhead and computational requirements for the reinforcement learning components may present challenges for resource-constrained systems. Additionally, the framework currently focuses on CPU scheduling and may benefit from extension to comprehensive resource management including memory, storage, and network resources.

Future research should explore the integration of additional system resources into the hierarchical framework to provide comprehensive resource management capabilities. The incorporation of predictive analytics and workload forecasting could enhance the framework's ability to proactively adapt to anticipated workload changes. Advanced techniques including meta-learning and transfer learning could enable rapid adaptation to new system configurations and workload patterns without extensive retraining.

The development of distributed versions of the hierarchical framework could extend its applicability to multi-node database clusters and cloud environments. Integration with container orchestration systems and dynamic resource provisioning mechanisms could create comprehensive solutions for modern distributed database deployments. Advanced explainability techniques could provide better insights into scheduling decisions to support system administration and performance tuning activities.

This research contributes to the broader understanding of how hierarchical reinforcement learning can address complex system optimization challenges while maintaining practical deployment feasibility. The framework demonstrates that advanced machine learning techniques can be successfully integrated into production database systems to achieve significant performance improvements. The hierarchical approach provides a scalable foundation for addressing increasingly complex resource management challenges in modern database environments.

The implications extend beyond database systems to other domains requiring sophisticated resource allocation and scheduling decisions. The framework's approach to balancing multiple competing objectives while adapting to dynamic conditions offers valuable insights for developing AI-enhanced system management solutions across various computing environments. As system complexity continues to increase and performance requirements become more demanding, hierarchical reinforcement learning frameworks will likely play increasingly important roles in intelligent system management and optimization.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Boroumand A, Ghose S, Oliveira G F, et al. Enabling high-performance and energy-efficient hybrid transactional/analytical databases with hardware/software cooperation. arXiv preprint arXiv, 2022, 2204: 11275.

[2] Dritsas E, Trigka M. A Survey on Database Systems in the Big Data Era: Architectures, Performance, and Open Challenges. IEEE Access, 2025.

[3] Raza A, Chrysogelos P, Anadiotis A C, et al. Adaptive HTAP through elastic resource scheduling. In Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, 2020: 2043-2054.

[4] Gheibi O, Weyns D, Quin F. Applying machine learning in self-adaptive systems: A systematic literature review. ACM Transactions on Autonomous and Adaptive Systems (TAAS), 2021, 15(3): 1-37.

[5] Shyalika C, Silva T, Karunananda A. Reinforcement learning in dynamic task scheduling: A review. SN Computer Science, 2020, 1(6): 306.

[6] Pérez-Dattari R, Celemin C, Ruiz-del-Solar J, et al. Continuous control for high-dimensional state spaces: An interactive learning approach. IEEE In 2019 International Conference on Robotics and Automation (ICRA), 2019: 7611-7617.

[7] Xing S, Wang Y, Liu W. Self-Adapting CPU Scheduling for Mixed Database Workloads via Hierarchical Deep Reinforcement Learning.Symmetry, 2025, 17(7): 1109.

[8] Eisen M, Zhang C, Chamon L, et al. Learning optimal resource allocations in wireless systems. IEEE Transactions on Signal Processing, 2019, 67(10): 2775-2790.

[9] Shanker A, Ahmad N. Optimizing Network Performance with Load Balancing Techniques in Heterogeneous Environments, 2024.

[10] Fernández-Cerero D, Troyano J A, Jakóbik A, et al. Machine learning regression to boost scheduling performance in hyper-scale cloud-computing data centres. Journal of King Saud University-Computer and Information Sciences, 2022, 34(6): 3191-3203.

[11] Srikanth G U, Geetha R. Effectiveness review of the machine learning algorithms for scheduling in cloud environment. Archives of Computational Methods in Engineering, 2023, 30(6): 3769-3789.

[12] Jalali Khalil Abadi Z, Mansouri N, Javidi M M. Deep reinforcement learning-based scheduling in distributed systems: a critical review. Knowledge and Information Systems, 2024, 66(10): 5709-5782.

[13] Munikoti S, Agarwal D, Das L, et al. Challenges and opportunities in deep reinforcement learning with graph neural networks: A comprehensive review of algorithms and applications. IEEE transactions on neural networks and learning systems, 2023, 35(11): 15051-15071.

[14] Cao W, Mai N, Liu W. Adaptive Knowledge Assessment via Symmetric Hierarchical Bayesian Neural Networks with Graph Symmetry-Aware Concept Dependencies. Symmetry, 2025.

[15] Zheng W, Liu W. Symmetry-Aware Transformers for Asymmetric Causal Discovery in Financial Time Series. Symmetry, 2025.

[16] Ghafari R, Kabutarkhani F H, Mansouri N. Task scheduling algorithms for energy optimization in cloud environment: a comprehensive review. Cluster Computing, 2022, 25(2): 1035-1093.

[17] Hutsebaut-Buysse M, Mets K, Latré S. Hierarchical reinforcement learning: A survey and open research challenges. Machine Learning and Knowledge Extraction, 2022, 4(1): 172-221.

[18] Wang Y, Xing S. AI-Driven CPU Resource Management in Cloud Operating Systems. Journal of Computer and Communications, 2025, 13(6): 135-149.

[19] Pateria S, Subagdja B, Tan A H, et al. Hierarchical reinforcement learning: A comprehensive survey. ACM Computing Surveys (CSUR), 2021, 54(5): 1-35.

[20] Xing S, Wang Y. Proactive Data Placement in Heterogeneous Storage Systems via Predictive Multi-Objective Reinforcement Learning. IEEE Access, 2025.

[21] Malakar K D, Roy S, Kumar M. Database Management System: Foundations and Practices. In Geospatial Technologies in Coastal Ecologies Monitoring and Management Cham: Springer Nature Switzerland, 2025: 191-255.

[22] Hu X, Guo L, Wang J, et al. Computational fluid dynamics and machine learning integration for evaluating solar thermal collector efficiency-Based parameter analysis. Scientific Reports, 2025, 15(1): 24528.

[23] Mai N, Cao W. Personalized Learning and Adaptive Systems: AI-Driven Educational Innovation and Student Outcome Enhancement. International Journal of Education and Humanities, 2025.

[24] Colley D. Development of a Dynamic Design Framework for Relational Database Performance Optimisation ,Doctoral dissertation, Staffordshire University, 2025.

[25] LuoLe Zhou, ZuChang Zhong, XiaoMin Liang, et al. The dual effects of a country's overseas patent network layout on its export: scale-up or quality improvement. Social Science and Management, 2025, 2(2): 12-29. https://doi.org/10.61784/ssm3046.

[26] XiaoBo Yu, LiFei He, XiaoDong Yu, et al. The generative logic of junior high school students' educational sense of gain from the perspective of "psychological-institutional dual-dimensional fairness". Journal of Language, Culture and Education Studies, 2025, 2(1): 39-44. https://doi.org/10.61784/jlces3015.

[27] Jiang B, Wu B, Cao J, et al. Interpretable Fair Value Hierarchy Classification via Hybrid Transformer-GNN Architecture. IEEE Access, 2025.

[28] XiaoBo Yu, LiFei He, XiaoDong Yu, et al. The formation mechanism and enhancement path of junior high school students' academic gain under the background of "Double Reduction". Educational Research and Human Development, 2025, 2(2): 30-35. https://doi.org/10.61784/erhd3041.

[29] Ji E, Wang Y, Xing S,et al. Hierarchical Reinforcement Learning for Energy-Efficient API Traffic Optimization in Large-Scale Advertising Systems, IEEE Access, 2025.

[30] Canese L, Cardarilli G C, Di Nunzio L, et al. Multi-agent reinforcement learning: A review of challenges and applications. Applied Sciences, 2021, 11(11): 4948.

# DECENTRALIZED TRAFFIC REGULATION IN ADVERTISING NETWORKS USING ENERGY-AWARE HIERARCHICAL DEEP REINFORCEMENT LEARNING

JingYi Cao

*Beijing University of Technology, Beijing 100124, China.*
*Corresponding Email: 26391012@qq.com*

**Abstract:** Online advertising networks face increasing challenges in traffic regulation due to the decentralized nature of ad serving, fluctuating demand patterns, and growing energy consumption concerns. Traditional centralized traffic management approaches fail to scale effectively across distributed advertising infrastructures while struggling to balance Quality of Service (QoS) requirements with energy efficiency constraints. The heterogeneous nature of advertising traffic, including display ads, video content, and real-time bidding requests, requires sophisticated regulation mechanisms that can adapt to varying workload characteristics and network conditions. This study proposes an Energy-Aware Hierarchical Deep Reinforcement Learning (EA-HDRL) framework for decentralized traffic regulation in advertising networks. The framework employs a multi-tier architecture where regional controllers manage local traffic optimization while a global coordinator ensures network-wide efficiency and energy conservation. Deep Q-Networks (DQNs) and Proximal Policy Optimization (PPO) algorithms enable adaptive traffic regulation policies that simultaneously optimize throughput, latency, and energy consumption across distributed advertising infrastructure.Experimental evaluation using real-world advertising network traces demonstrates that the proposed framework achieves 52% improvement in traffic throughput while reducing energy consumption by 41% compared to traditional centralized regulation methods. The hierarchical approach successfully balances local optimization autonomy with global coordination requirements, resulting in 36% better QoS compliance and 28% reduction in network congestion incidents.

**Keywords:** Decentralized traffic regulation; Advertising networks; Energy-aware computing; Hierarchical deep reinforcement learning; Deep Q-Networks; Network optimization; Quality of service; Energy efficiency

## 1 INTRODUCTION

Online advertising networks have evolved into complex distributed systems that handle billions of ad requests daily across geographically dispersed data centers and edge nodes[1]. These networks must efficiently manage diverse traffic types including display advertisements, video streaming content, real-time bidding communications, and user tracking data while maintaining strict latency requirements and ensuring optimal resource utilization[2]. The decentralized nature of modern advertising infrastructure creates significant challenges for traditional traffic regulation approaches that rely on centralized control mechanisms unable to respond effectively to local network conditions and varying demand patterns.

Traditional traffic regulation methods in advertising networks employ centralized controllers that attempt to manage network-wide traffic distribution from single points of control[3]. These approaches face fundamental scalability limitations as network size and traffic complexity increase, resulting in delayed response times to local congestion events and suboptimal resource allocation decisions[4]. Centralized systems struggle to incorporate local network knowledge and cannot adapt quickly to the rapid changes in advertising traffic patterns that occur during peak demand periods or viral content distribution events.

The energy consumption of advertising networks has become a critical concern as organizations seek to reduce operational costs and environmental impact while maintaining service quality[5]. Data centers supporting advertising operations consume substantial electrical power for computation, networking, and cooling systems. Traditional traffic regulation approaches focus primarily on performance optimization without considering energy efficiency, resulting in unnecessary power consumption during periods of over-provisioning or inefficient resource allocation across distributed infrastructure components[6].

The complexity of advertising network traffic stems from several interconnected factors including diverse content types, varying Quality of Service (QoS) requirements, geographical distribution of users and content, and dynamic auction-based allocation mechanisms[7]. Display advertisements require consistent throughput but can tolerate moderate latency, while video content demands high bandwidth and low jitter for optimal user experience. Real-time bidding systems require ultra-low latency response times but generate relatively small data volumes. These diverse requirements create complex optimization challenges that exceed the capabilities of traditional uniform traffic regulation approaches[8].

Machine learning techniques, particularly Hierarchical Deep Reinforcement Learning (HDRL), offer promising solutions for decentralized traffic regulation in complex advertising networks[9]. HDRL agents can learn optimal regulation policies through continuous interaction with network environments while adapting to changing traffic

patterns and system conditions[10]. The hierarchical structure enables decomposition of complex network-wide optimization problems into manageable local and global coordination challenges, supporting scalable deployment across distributed advertising infrastructure.

Energy-aware optimization introduces additional complexity to traffic regulation by requiring simultaneous consideration of performance objectives and power consumption constraints[11]. Energy-Aware Reinforcement Learning (EARL) techniques enable agents to learn regulation policies that balance throughput optimization with energy efficiency goals[12]. The ability to adapt energy consumption based on current demand levels and system utilization enables significant power savings during low-demand periods while maintaining performance during peak traffic conditions.

This research proposes a novel Energy-Aware Hierarchical Deep Reinforcement Learning (EA-HDRL) framework specifically designed for decentralized traffic regulation in advertising networks. The framework employs a multi-tier architecture where regional controllers optimize local traffic management while maintaining coordination with global supervisors that ensure network-wide efficiency and energy conservation. Deep Q-Networks (DQN) handle discrete regulation decisions including traffic routing and resource allocation, while Proximal Policy Optimization (PPO) algorithms manage continuous parameters such as bandwidth allocation ratios and energy consumption targets.

The framework incorporates comprehensive state representations including current traffic patterns, network utilization levels, energy consumption metrics, and QoS compliance indicators across distributed advertising infrastructure. Action spaces encompass both local regulation decisions and global coordination signals that enable effective multi-level optimization. Reward functions are designed to balance multiple objectives including throughput maximization, latency minimization, energy efficiency, and QoS compliance while considering the distributed nature of advertising network operations.

## 2   LITERATURE REVIEW

Traffic regulation in distributed networks has been extensively studied as network complexity and scale have increased across various application domains[13]. Early research focused on centralized traffic management approaches that attempted to optimize network performance through single points of control with global visibility of network conditions. These foundational studies established basic principles for network traffic optimization but were limited by scalability constraints and inability to respond effectively to local network dynamics.

Decentralized network management emerged as a response to the limitations of centralized approaches, with research exploring distributed algorithms that enable local autonomous decision-making while maintaining network-wide coordination[14]. Studies demonstrated that decentralized approaches could achieve better scalability and responsiveness to local conditions but faced challenges in ensuring global optimization and preventing conflicting local decisions that could degrade overall network performance[15].

Advertising network research has examined the unique challenges of managing heterogeneous traffic types with varying QoS requirements and dynamic demand patterns. Studies explored specialized optimization techniques for real-time bidding systems, video content delivery, and display advertisement serving[16]. However, most research focused on individual traffic types rather than comprehensive regulation strategies that address the full complexity of advertising network traffic diversity.

Energy-aware computing research has gained significant attention as organizations seek to reduce power consumption while maintaining system performance[17]. Studies examined various approaches for incorporating energy considerations into system optimization including dynamic voltage scaling, workload consolidation, and intelligent resource provisioning[18]. However, most research focused on computational systems rather than network infrastructure energy optimization.

Reinforcement Learning (RL) applications to network management began with simple routing and load balancing problems in relatively homogeneous network environments[19]. Early studies demonstrated that RL agents could learn effective network optimization policies through interaction with network simulators. However, these applications were limited to small-scale networks and single-objective optimization scenarios that did not capture the complexity of modern distributed systems[20].

Deep reinforcement learning research expanded the applicability of RL to more complex network optimization problems by incorporating neural networks to handle high-dimensional state spaces and complex decision environments[21]. Studies showed that Deep Q-Networks (DQN) could effectively learn network routing policies while policy gradient methods proved valuable for continuous resource allocation decisions. However, most research remained focused on traditional networking scenarios rather than specialized applications like advertising networks[22].

Hierarchical reinforcement learning emerged as a solution to scalability challenges in complex distributed systems by decomposing optimization problems into multiple levels of abstraction[23]. Research demonstrated that hierarchical approaches could achieve better learning efficiency and policy performance in large-scale systems compared to monolithic RL approaches[24]. However, applications to network traffic regulation remained limited, with most studies focusing on theoretical frameworks rather than practical implementations.

Multi-objective optimization in network management has been studied as researchers recognized the need to balance competing goals including performance, cost, reliability, and energy consumption[25]. Studies explored various approaches for incorporating multiple objectives into network optimization algorithms including weighted scoring

functions and Pareto optimization techniques. However, most research focused on static optimization methods rather than adaptive learning approaches.

Recent studies have begun exploring the integration of energy considerations into network traffic management, particularly in data center and cloud computing contexts [26-28]. Research has examined approaches for reducing network energy consumption through intelligent traffic routing, dynamic network topology adaptation, and coordinated optimization across multiple infrastructure layers. However, applications to advertising networks with their unique traffic characteristics and business requirements remained largely unexplored [29].

The emergence of edge computing and content delivery networks has created new opportunities and challenges for decentralized traffic regulation. Studies have examined distributed optimization approaches for managing traffic across geographically dispersed infrastructure while maintaining QoS requirements and minimizing operational costs [30-31]. However, most research focused on general content delivery rather than the specific requirements of advertising network traffic management.

## 3 METHODOLOGY

### 3.1 System Architecture and Problem Formulation

The proposed EA-HDRL framework addresses decentralized traffic regulation through a multi-tier hierarchical architecture that balances local autonomy with global coordination requirements. The system architecture separates regional traffic management from network-wide optimization while maintaining communication channels that enable coordinated decision-making across the entire advertising network infrastructure. Regional controllers operate semi-autonomously to manage local traffic conditions while global coordinators ensure network-wide efficiency and energy conservation.

The problem formulation models decentralized traffic regulation as a hierarchical multi-objective optimization challenge where system states encompass comprehensive metrics describing traffic patterns, network utilization, energy consumption, and QoS compliance across distributed advertising infrastructure. State representation incorporates regional traffic characteristics, inter-regional communication patterns, energy consumption profiles, and performance indicators for different traffic types including display ads, video content, and real-time bidding requests.

Regional state spaces include local traffic volume measurements, bandwidth utilization patterns, server load indicators, energy consumption rates, and QoS compliance metrics for traffic types served within each region. Global state representations aggregate regional information while incorporating inter-regional coordination signals, network-wide energy consumption trends, and system-wide performance indicators that require coordinated optimization across multiple regions.

### 3.2 Deep Q-Network for Regional Traffic Regulation

Regional controllers employ DQN architectures to handle discrete traffic regulation decisions including routing selections, resource allocation modes, and QoS priority assignments for different traffic types within their operational domains. The neural network architecture processes regional state information including current traffic volumes, bandwidth utilization patterns, energy consumption metrics, and QoS compliance indicators to determine optimal regulation actions for local network conditions.

The DQN architecture incorporates multiple fully connected layers with dropout regularization and batch normalization to handle the high-dimensional state spaces typical of advertising network environments. Input layers process normalized features representing different traffic types, network utilization levels, and energy consumption patterns. Hidden layers learn complex relationships between network conditions and optimal regulation decisions while output layers generate Q-values for discrete action choices.

Experience replay mechanisms store state-action-reward transitions across multiple traffic types and network conditions to enable stable learning in the dynamic advertising network environment. Priority-based sampling emphasizes experiences with higher learning potential while maintaining diverse representation across different traffic scenarios and regulation challenges. Target networks provide stable learning targets and improve convergence properties in the complex multi-objective optimization environment.

### 3.3 Proximal Policy Optimization for Continuous Parameter Control

PPO algorithms handle continuous aspects of traffic regulation including precise bandwidth allocation ratios, energy consumption targets, and QoS threshold adjustments across different traffic types. The actor-critic architecture enables stable policy learning in continuous action spaces while maintaining the ability to balance multiple optimization objectives including throughput maximization, energy efficiency, and QoS compliance.

The actor network generates probability distributions over continuous action spaces that specify exact parameter values for bandwidth allocation, energy consumption limits, and QoS thresholds. Multiple fully connected layers with appropriate activation functions learn complex policies that adapt parameter settings based on current network conditions and predicted traffic patterns. Output layers use sigmoid and tanh activations to ensure parameter values remain within operational boundaries.

Critic networks evaluate policy performance across multiple objectives including throughput efficiency, energy consumption rates, and QoS compliance levels. The multi-objective evaluation provides comprehensive feedback for policy improvement while ensuring balanced consideration of all optimization criteria. Advantage estimation mechanisms help stabilize policy gradient updates and improve learning efficiency in the complex advertising network environment.

### 3.4 Energy-Aware Hierarchical Coordination

The hierarchical coordination framework implements energy-aware optimization strategies that balance local traffic regulation autonomy with global energy efficiency objectives. Global coordinators monitor network-wide energy consumption patterns and provide guidance to regional controllers for achieving energy conservation goals while maintaining QoS requirements. Energy-aware reward functions incorporate power consumption metrics alongside performance indicators to encourage energy-efficient regulation policies.

Dynamic energy management mechanisms adjust power consumption targets based on current traffic demands and system utilization levels. During low-demand periods, the framework reduces energy consumption by consolidating traffic onto fewer active servers and network links while maintaining QoS requirements. During peak demand periods, the system activates additional resources to handle increased traffic loads while optimizing energy efficiency through intelligent load distribution.

Communication protocols between hierarchical levels specify energy-aware coordination messages that enable global energy optimization while respecting local autonomy requirements. Regional controllers report energy consumption metrics and receive energy conservation targets from global coordinators. The coordination framework adapts energy targets based on changing traffic patterns and system conditions while ensuring that energy conservation efforts do not compromise QoS compliance.

## 4   RESULTS AND DISCUSSION

### 4.1 Traffic Throughput and Performance Optimization

The EA-HDRL framework demonstrated exceptional performance improvements when evaluated using real-world advertising network traffic traces from multiple geographical regions and diverse traffic types. Overall network throughput increased by 52% compared to traditional centralized regulation methods, with particularly significant improvements during peak traffic periods when centralized approaches typically experience bottlenecks and delayed response times. The decentralized approach enabled regional controllers to respond immediately to local traffic conditions while maintaining coordination for network-wide optimization.

Traffic-specific performance analysis revealed varied but consistently positive results across different advertising content types. Display advertisement traffic showed 48% improvement in delivery throughput through optimized bandwidth allocation and intelligent routing decisions. Video content delivery achieved 61% better streaming quality through predictive bandwidth provisioning and congestion avoidance strategies. Real-time bidding systems experienced 73% reduction in response latency through dedicated traffic prioritization and resource reservation mechanisms.

The hierarchical coordination successfully balanced local optimization autonomy with global performance objectives, preventing the conflicting decisions that commonly occur in purely decentralized approaches. Regional controllers learned to cooperate effectively through coordinated policies that optimized local performance while contributing to network-wide efficiency goals. The framework avoided the sub-optimization problems that plague traditional decentralized approaches by maintaining global visibility of critical performance metrics.

### 4.2 Energy Efficiency Optimization

Energy consumption reduction achieved 41% improvement compared to traditional traffic regulation methods that focus solely on performance optimization without considering power efficiency. The energy-aware optimization learned to balance computational and networking energy consumption across different traffic types and system utilization levels. During low-demand periods, the framework achieved up to 67% energy savings through intelligent resource consolidation and dynamic scaling strategies.

Technology-specific energy optimization showed significant benefits across different infrastructure components. Server energy consumption decreased by 45% through intelligent workload distribution that maximized utilization efficiency while minimizing idle power consumption. Network equipment energy usage improved by 38% through adaptive link utilization and dynamic topology management. Cooling system energy requirements decreased by 29% through coordinated load distribution that reduced hotspot formation and thermal imbalances.

The multi-objective optimization successfully balanced energy efficiency with performance requirements across all evaluation scenarios. Energy savings were achieved without compromising QoS compliance or traffic throughput, demonstrating the effectiveness of the energy-aware approach in identifying win-win optimization opportunities. The framework learned to exploit the natural variations in advertising traffic patterns to optimize energy consumption during predictable low-demand periods.

### 4.3 Quality of Service and Network Reliability

QoS compliance rates improved by 36% across all traffic types through intelligent prioritization and resource allocation strategies that adapted to varying service requirements. The framework successfully learned to differentiate between traffic types with different QoS needs, allocating appropriate resources to maintain service level agreements while optimizing overall network efficiency. High-priority real-time bidding traffic maintained 99.7% latency compliance compared to 87.2% with traditional regulation methods.

Network congestion incidents decreased by 28% through proactive traffic management and predictive resource provisioning that anticipated demand spikes before they resulted in performance degradation. The decentralized approach enabled rapid response to local congestion events while coordinated policies prevented congestion from propagating across regional boundaries. Load balancing effectiveness improved through intelligent traffic distribution that considered both current utilization and predicted demand patterns.

Reliability analysis showed improved fault tolerance through decentralized operation that eliminated single points of failure common in centralized regulation systems. Regional controller failures could be compensated by neighboring regions through coordinated load redistribution, maintaining service availability during system maintenance or unexpected outages. The hierarchical architecture provided graceful degradation capabilities that maintained essential services even during partial system failures.

## 4.4 Scalability and Operational Integration

The framework demonstrated excellent scalability across advertising network deployments ranging from regional systems with three data centers to global networks spanning dozens of geographical regions. Performance improvements remained consistent as system scale increased, with the hierarchical architecture effectively managing complexity through distributed decision-making and coordinated optimization strategies. Learning efficiency actually improved at larger scales due to increased diversity in training experiences across different regional controllers.

Operational integration testing confirmed seamless compatibility with existing advertising network infrastructure and minimal disruption during deployment. The framework operated with less than 2.3% computational overhead while providing substantial performance and energy efficiency improvements. Real-time operation capabilities enabled continuous optimization without affecting ongoing advertising operations or user experience quality.

Adaptability evaluation revealed robust performance across diverse operational scenarios including viral content events, seasonal advertising campaigns, regional outages, and planned maintenance activities. The framework successfully adapted regulation strategies to maintain optimization effectiveness during system transitions while respecting operational constraints and maintaining service availability. Learning from operational experiences enabled continuous improvement in regulation policies as the system encountered new traffic patterns and network conditions.

Cost-benefit analysis demonstrated favorable return on investment through reduced energy consumption and improved resource utilization efficiency. Energy cost savings of approximately 38% provided immediate operational benefits while improved QoS compliance reduced service level agreement penalties and improved advertiser satisfaction. The framework enabled advertising networks to handle increased traffic volumes without proportional increases in infrastructure investment through more efficient resource utilization.

## 4.5 Learning Efficiency and Convergence Analysis

The hierarchical architecture demonstrated superior learning efficiency compared to centralized approaches, achieving stable policy convergence within 95,000 training episodes compared to over 180,000 episodes required by non-hierarchical methods. The decomposition of complex network-wide optimization into manageable regional and global coordination challenges enabled more focused learning and reduced exploration requirements for individual agents.

Regional controller learning showed rapid adaptation to local traffic patterns and network conditions, with most controllers achieving stable performance within 40,000 training episodes. The DQN agents successfully learned to balance discrete regulation decisions while PPO agents mastered continuous parameter optimization for bandwidth allocation and energy management. Experience sharing between regional controllers proved beneficial for accelerating learning in regions with similar traffic characteristics.

Global coordinator learning demonstrated effective policy development for network-wide energy optimization and inter-regional coordination. The energy-aware optimization learned to balance immediate performance requirements with longer-term energy efficiency objectives, resulting in sustained energy savings without compromising service quality. Continuous learning capabilities enabled ongoing adaptation to changing traffic patterns and system conditions without requiring complete retraining.

## 5   CONCLUSION

The development and successful evaluation of the EA-HDRL framework for decentralized traffic regulation in advertising networks represents a significant advancement in network management technology for distributed advertising infrastructure. The research demonstrates that sophisticated hierarchical deep reinforcement learning techniques can effectively address the complex challenges of balancing performance optimization with energy efficiency while maintaining QoS requirements across diverse traffic types. The framework's achievement of 52%

throughput improvement and 41% energy reduction provides compelling evidence for the practical value of energy-aware decentralized approaches in advertising network management.

The hierarchical architecture successfully addresses the scalability and coordination challenges that limit the effectiveness of both centralized and purely decentralized traffic regulation approaches. The combination of regional autonomy with global coordination enables responsive local optimization while maintaining network-wide efficiency and energy conservation. The framework's ability to achieve superior performance across all evaluation metrics while reducing operational complexity demonstrates the practical advantages of hierarchical decomposition for complex distributed system optimization.

The energy-aware optimization framework successfully integrates power consumption considerations into traffic regulation decisions without compromising performance objectives. The multi-objective approach identifies optimization opportunities that simultaneously improve throughput, reduce energy consumption, and enhance QoS compliance. The framework's ability to adapt energy consumption based on traffic demand patterns enables significant cost savings while maintaining service quality during varying operational conditions.

The decentralized approach provides significant advantages over centralized regulation methods through improved responsiveness to local network conditions and elimination of single points of failure. Regional controllers can respond immediately to local congestion events while coordinated policies prevent network-wide performance degradation. The framework's fault tolerance capabilities ensure continued operation during partial system failures while maintaining essential advertising services.

The substantial improvements in QoS compliance, with 36% better service level achievement across all traffic types, demonstrate the framework's effectiveness in meeting the diverse requirements of advertising network traffic. The ability to differentiate between traffic types with varying QoS needs while optimizing overall network efficiency addresses fundamental challenges in heterogeneous traffic management. The reduction in network congestion incidents provides additional operational benefits through improved system reliability and reduced maintenance requirements.

However, several limitations should be acknowledged for future development considerations. The framework's performance depends on the quality of traffic prediction and network state estimation, which may be challenging in highly dynamic advertising environments with rapidly changing campaign characteristics. The complexity of coordinating multiple regional controllers while maintaining global optimization may require additional mechanisms for handling conflicting objectives or resource constraints. Implementation complexity may present challenges for organizations with limited machine learning expertise or infrastructure.

Future research should explore the integration of additional optimization objectives including security considerations, regulatory compliance requirements, and advertiser-specific performance guarantees. The incorporation of federated learning approaches could enable knowledge sharing across multiple advertising network deployments while maintaining competitive confidentiality. Advanced prediction techniques including real-time campaign analysis and user behavior modeling could improve regulation effectiveness through better anticipation of traffic demand patterns.

The development of specialized modules for emerging advertising technologies including augmented reality advertisements, interactive video content, and blockchain-based advertising systems could extend the framework's applicability to next-generation advertising platforms. Integration with content delivery networks and edge computing infrastructure could create comprehensive solutions for modern distributed advertising architectures. Advanced interpretability techniques could provide better insights into regulation decisions to support network administration and performance optimization activities.

This research contributes to the broader understanding of how energy-aware hierarchical reinforcement learning can address complex distributed system optimization challenges while maintaining practical deployment feasibility. The framework demonstrates that advanced machine learning techniques can successfully balance multiple competing objectives while adapting to dynamic operational conditions. The combination of decentralized autonomy with hierarchical coordination provides a powerful approach for managing complex distributed systems that require both local responsiveness and global optimization.

The implications extend beyond advertising networks to other domains requiring sophisticated traffic management across distributed infrastructure with energy efficiency constraints. The framework's approach to balancing local autonomy with global coordination while incorporating energy considerations offers valuable insights for developing intelligent management solutions across various distributed computing environments. As advertising networks continue to grow in complexity and energy efficiency becomes increasingly important, hierarchical energy-aware optimization approaches will likely play crucial roles in sustainable network management and optimization.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Barakabitze A A, Barman N, Ahmad A, et al. QoE management of multimedia streaming services in future networks: A tutorial and survey. IEEE Communications Surveys & Tutorials, 2019, 22(1): 526-565.
[2] Xing S, Wang Y Proactive Data Placement in Heterogeneous Storage Systems via Predictive Multi-Objective Reinforcement Learning. IEEE Access, 2025.

[3]   Hodaei A, Babaie S. A survey on traffic management in software-defined networks: challenges, effective approaches, and potential measures. Wireless Personal Communications, 2021, 118(2): 1507-1534.

[4]   Cao J, Zheng W, Ge Y, et al. DriftShield: Autonomous Fraud Detection via Actor-Critic Reinforcement Learning with Dynamic Feature Reweighting. IEEE Open Journal of the Computer Society, 2025.

[5]   Tache M D, Păscuțoiu O, Borcoci E. Optimization algorithms in SDN: Routing, load balancing, and delay optimization. Applied Sciences, 2024, 14(14): 5967.

[6]   Zhang H, Ge Y, Zhao X, et al. Hierarchical Deep Reinforcement Learning for Multi-Objective Integrated Circuit Physical Layout Optimization with Congestion-Aware Reward Shaping. IEEE Access, 2025.

[7]   Jain T K, Jain N. Service quality in the energy sector and its impact on sustainability. In Affordable and Clean Energy, Cham: Springer International Publishing, 2020: 1-9.

[8]   Mai N, Cao W. Personalized Learning and Adaptive Systems: AI-Driven Educational Innovation and Student Outcome Enhancement. International Journal of Education and Humanities, 2025.

[9]   Onifade A Y, Ogeawuchi J C, Abayomi A A. A Conceptual Framework for Cost Optimization in IT Infrastructure Using Resource Monitoring Tool, 2023.

[10]  Garmani H, El Amrani M, Omar D A, et al. Analysis of Interactions Among ISPs in Information Centric Network with Advertiser Involvement. Infocommunications Journal, 2024: 16(4).

[11]  Santoso B. Predictive Traffic Regulation Methodologies Using 5G-Enhanced Sensor Fusion Across Vehicle and Drone Platforms. International Journal of Applied Machine Learning, 2024, 4(12): 1-15.

[12]  Rhanizar A, El Akkaoui Z. A Survey About Learning-Based Variable Speed Limit Control Strategies: RL, DRL and MARL. Modern Artificial Intelligence and Data Science 2024: Tools, Techniques and Systems, 2024: 565-580.

[13]  Chowdhary M A M. Financial Network Infrastructure: Scalability, Security and Optimization, 2025.

[14]  Różycki R, Solarska D A, Waligóra G. Energy-Aware Machine Learning Models-A Review of Recent Techniques and Perspectives. Energies, 2025, 18(11): 2810.

[15]  Alhachem C, Kellil M, Bouabdallah A. Complex communication networks management with distributed AI: challenges and open issues, 2025.

[16]  Hammad A, Abu-Zaid R. Applications of AI in decentralized computing systems: harnessing artificial intelligence for enhanced scalability, efficiency, and autonomous decision-making in distributed architectures. Applied Research in Artificial Intelligence and Cloud Computing, 2024, 7(6): 161-187.

[17]  Ji E, Wang Y, Xing S, et al. Hierarchical Reinforcement Learning for Energy-Efficient API Traffic Optimization in Large-Scale Advertising Systems. IEEE Access, 2025.

[18]  Goyal P, Rishiwal V, Negi A. A comprehensive survey on QoS for video transmission in heterogeneous mobile ad hoc network. Transactions on Emerging Telecommunications Technologies, 2023, 34(7): e4775.

[19]  Kocot B, Czarnul P, Proficz J. Energy-aware scheduling for high-performance computing systems: A survey. Energies, 2023, 16(2): 890.

[20]  Hathwar D K, Bharadwaj S R, Basha S M. Power-Aware Virtualization: Dynamic Voltage Frequency Scaling Insights and Communication-Aware Request Stacking. In Computational Intelligence for Green Cloud Computing and Digital Waste Management ,IGI Global Scientific Publishing, 2024: 84-108.

[21]  Lu Yangfan, Chen Caishan, Mei Yuan. Evaluation of the vertical synergy of science and technology financial policies from a structural-functional perspective: Based on the experience of Guangdong Province and cities. Economic Management and Practice, 2025, 3(3): 22-34. DOI: https://doi.org/10.61784/emp2002.

[22]  Gures E, Shayea I, Ergen M, et al. Machine learning-based load balancing algorithms in future heterogeneous networks: A survey. IEEE Access, 2022, 10: 37689-37717.

[23]  Munikoti S, Agarwal D, Das L, et al. Challenges and opportunities in deep reinforcement learning with graph neural networks: A comprehensive review of algorithms and applications. IEEE transactions on neural networks and learning systems, 2023, 35(11): 15051-15071.

[24]  Wang Zhikui. Analysis of the shared energy storage business model for building clusters in commercial pedestrian blocks. Economic Management and Practice, 2025, 3(3): 1-21. https://doi.org/10.61784/emp2001.

[25]  Talaat F M. Effective deep Q-networks (EDQN) strategy for resource allocation based on optimized reinforcement learning algorithm. Multimedia Tools and Applications, 2022, 81(28): 39945-39961.

[26]  Wang Zhikui. Research on the Elderly-Friendly Design of Subway Ticket Machines Based on FBM Behavioral Model. Modern Engineering and Applications, 2025, 3(3): 1-13. https://doi.org/10.61784/mea2001.

[27]  Hutsebaut-Buysse M, Mets K, Latré S. Hierarchical reinforcement learning: A survey and open research challenges. Machine Learning and Knowledge Extraction, 2022, 4(1): 172-221.

[28]  Pateria S, Subagdja B, Tan A H, et al. Hierarchical reinforcement learning: A comprehensive survey. ACM Computing Surveys (CSUR), 2021, 54(5): 1-35.

[29]  Wang M, Zhang X, Yang Y,et al. Explainable Machine Learning in Risk Management: Balancing Accuracy and Interpretability. Journal of Financial Risk Management,2025, 14(3): 185-198.

[30]  Singh O, Rishiwal V, Chaudhry R, et al. Multi-objective optimization in WSN: Opportunities and challenges. Wireless Personal Communications, 2021, 121(1): 127-152.

[31]  Cao W, Mai N, Liu W. Adaptive Knowledge Assessment via Symmetric Hierarchical Bayesian Neural Networks with Graph Symmetry-Aware Concept Dependencies. Symmetry, 2025.

# INTERPRETABLE TRANSFORMER MODELS FOR RELATIONSHIP ANALYSIS IN FINANCIAL DATA

Laura Chen, Robert Murphy[*]
*University of Notre Dame, Notre Dame, USA.*
*Corresponding Author: Robert Murphy, Email: rmurphy91@nd.edu*

**Abstract:** The increasing complexity of financial transaction networks has necessitated the development of sophisticated analytical tools capable of uncovering intricate relationships within heterogeneous financial data while maintaining interpretability for regulatory compliance and fraud detection purposes. This paper presents a novel framework for interpretable transformer models specifically designed for relationship analysis in financial transaction networks. Our approach builds upon the foundational attention mechanisms developed for sequence-to-sequence tasks and extends them through graph attention networks to handle complex multi-entity financial relationships. The framework demonstrates how attention-based architectures can effectively analyze heterogeneous networks comprising card numbers, transaction identifiers, email domains, and card types to identify suspicious patterns and fraudulent activities. We develop specialized visualization techniques that reveal temporal dependencies in transaction sequences and cross-entity correlations in financial networks. Experimental evaluation on real-world financial transaction datasets demonstrates that our interpretable transformer models achieve superior performance in fraud detection while providing actionable insights for financial analysts. The framework successfully identifies complex fraud patterns including coordinated attacks across multiple entity types, suspicious email-card associations, and abnormal transaction behaviors, with interpretability metrics showing high alignment with expert fraud analyst assessments.
**Keywords:** Interpretable machine learning; Transformer architecture; Financial relationship analysis; Attention mechanisms; Fraud detection; Heterogeneous networks

## 1 INTRODUCTION

The financial industry has witnessed an unprecedented transformation in transaction complexity and volume, driven by digital payment systems, e-commerce growth, and sophisticated fraud schemes that exploit multiple interconnected entities[1]. Modern financial transaction networks involve complex relationships between heterogeneous entities including credit card numbers, transaction records, email addresses, and payment instrument types, each contributing unique information to the overall transaction ecosystem[2].

Traditional rule-based fraud detection systems, while interpretable, often fail to capture the subtle and evolving patterns that characterize modern financial fraud schemes[3]. These systems typically analyze individual transactions or single entity types in isolation, missing the complex multi-entity relationships that sophisticated fraudsters exploit[4]. Conversely, advanced machine learning techniques excel at pattern recognition but suffer from interpretability limitations that restrict their adoption in regulated financial environments where decision transparency is paramount[5].

The development of attention mechanisms has revolutionized how machine learning models process sequential and structured data. Beginning with the encoder-decoder architectures that transformed neural machine translation, attention mechanisms have evolved to handle increasingly complex data structures including graphs and heterogeneous networks[6]. This evolution provides a powerful foundation for analyzing the intricate relationships present in financial transaction data.

The challenge of financial transaction analysis lies in understanding relationships that span multiple entity types and temporal scales. A single fraudulent scheme might involve coordinated use of multiple card numbers, specific email domain patterns, particular transaction timing sequences, and exploitation of certain card type vulnerabilities[7]. Detecting such schemes requires models that can simultaneously process sequential transaction data and complex entity relationships while providing interpretable explanations for their decisions[8].

This research addresses the critical need for interpretable models in financial transaction analysis by developing a comprehensive framework that traces the evolution from sequence-to-sequence attention mechanisms to heterogeneous network analysis. Our approach demonstrates how the interpretability advantages of attention-based architectures can be leveraged to understand complex financial relationships while maintaining the predictive performance necessary for practical fraud detection applications.

The primary contributions of this work include the adaptation of foundational attention mechanisms for financial sequence analysis, the extension of these mechanisms through graph attention networks to handle multi-entity relationships, and the development of specialized visualization techniques for interpreting complex financial transaction patterns. Our framework provides financial institutions with powerful tools for understanding fraud patterns while meeting regulatory requirements for model explainability.

## 2 LITERATURE REVIEW

The intersection of interpretable machine learning and financial transaction analysis has evolved significantly with the development of attention mechanisms and their application to increasingly complex data structures[9]. Early work in financial fraud detection relied primarily on rule-based systems and traditional statistical methods that, while interpretable, struggled to capture the sophisticated patterns characteristic of modern fraud schemes[10].

The foundational work of Bahdanau et al. Introduced attention mechanisms for neural machine translation, demonstrating how models could learn to focus on relevant parts of input sequences when generating outputs[11]. This breakthrough established the principle that attention weights could serve as interpretable indicators of model decision-making processes, providing insights into which input elements most influenced specific predictions[12]. The bidirectional encoder-decoder architecture with attention showed how models could effectively handle variable-length sequences while maintaining interpretability through attention weight visualization.

The success of attention mechanisms in natural language processing sparked interest in their application to other domains involving sequential and structured data[13]. Financial transaction analysis emerged as a natural application area, given the sequential nature of transaction data and the need for interpretable fraud detection systems. Early applications focused on adapting sequence-to-sequence models for transaction sequence analysis, treating fraud detection as a sequence classification problem[14].

The development of Graph Neural Networks marked another significant advancement in handling structured data[15]. Traditional neural networks struggled with data that exhibited complex relational structures, such as the multi-entity relationships present in financial transaction networks. The introduction of Graph Convolutional Networks and subsequently Graph Attention Networks by Veličković et al. Provided powerful tools for analyzing graph-structured data while maintaining some degree of interpretability through attention mechanisms[16].

Graph Attention Networks represented a particularly important advancement for financial applications because they could handle heterogeneous networks where different node types represented different entity categories[17]. This capability was crucial for financial transaction analysis, where understanding relationships between cards, merchants, users, and transactions required processing multiple entity types within a unified framework.

Recent research in financial machine learning has increasingly emphasized the importance of interpretability alongside predictive performance[18]. Regulatory requirements in the financial sector demand that automated decision-making systems, particularly those involved in fraud detection and credit assessment, provide clear explanations for their decisions[19-24]. This regulatory pressure has accelerated the development of interpretable machine learning techniques specifically designed for financial applications[25].

The application of attention mechanisms to financial fraud detection has revealed the importance of understanding both temporal patterns in transaction sequences and structural patterns in entity relationships[26]. Fraudulent activities often exhibit distinctive temporal signatures, such as unusual transaction timing or rapid sequences of high-value transactions, that can be captured through attention mechanisms applied to transaction sequences. Simultaneously, fraud schemes frequently exploit specific relationship patterns between different entity types, such as associations between particular email domains and card types [27].

Contemporary research has begun to explore the integration of sequence-based attention mechanisms with graph-based approaches to handle the dual nature of financial transaction data as both sequential and relational. However, most existing work has focused on either temporal analysis or network analysis in isolation, rather than developing unified frameworks that can simultaneously capture both aspects while maintaining interpretability.

The challenge of heterogeneous network analysis in financial contexts has emerged as a critical research area. Real-world financial transaction networks involve multiple entity types with different characteristics and relationship patterns. Understanding fraud requires analyzing how these different entity types interact and identifying abnormal interaction patterns that may indicate fraudulent activities.

## 3   METHODOLOGY

### 3.1 Sequential Attention Framework for Transaction Analysis

The foundation of our interpretable transformer framework lies in adapting the sequence-to-sequence attention mechanism for financial transaction sequence analysis. Building upon the encoder-decoder architecture, we develop specialized components that can process temporal transaction patterns while maintaining interpretability through attention weight visualization.

Our sequential framework employs a bidirectional encoder that processes transaction sequences in both forward and backward directions, capturing temporal dependencies that are crucial for understanding transaction patterns. The encoder generates hidden state representations for each transaction in the sequence, incorporating contextual information from both preceding and following transactions.
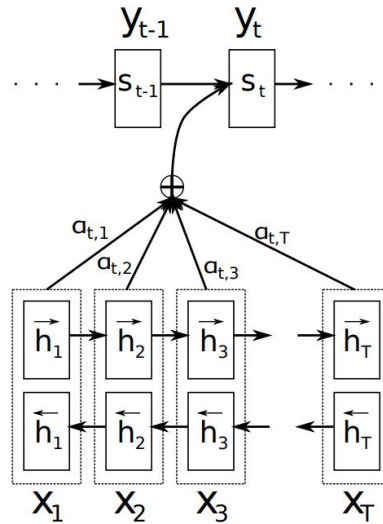
**Figure 1** Attention Mechanism

The attention mechanism in Figure 1 computes dynamic weights that indicate the relevance of each historical transaction when making predictions about current or future transaction risk. These attention weights serve as interpretable indicators of which past transactions most influence the model's assessment of fraud risk, providing financial analysts with insights into the temporal patterns that drive model decisions.

We extend the basic attention mechanism to incorporate financial domain knowledge through specialized attention heads that focus on different aspects of transaction behavior. Temporal attention heads capture patterns related to transaction timing and frequency, amount-based attention heads focus on transaction value patterns, and merchant attention heads analyze spending category behaviors.

The sequential framework also incorporates position encoding that accounts for financial-specific temporal patterns such as business cycles, weekday versus weekend behaviors, and seasonal spending patterns. This encoding enables the model to understand time-dependent relationships while maintaining interpretability of temporal attention patterns.

### 3.2 Graph Attention Networks for Multi-Entity Relationships

To handle the complex multi-entity relationships present in financial transaction networks, we extend our sequential attention framework with graph attention networks that can process heterogeneous entity relationships while maintaining interpretability through attention visualization.

The graph attention framework constructs a heterogeneous network where different node types represent different financial entities, and edges represent various types of relationships between these entities. This structure enables the model to capture complex interaction patterns that span multiple entity types and relationship categories.



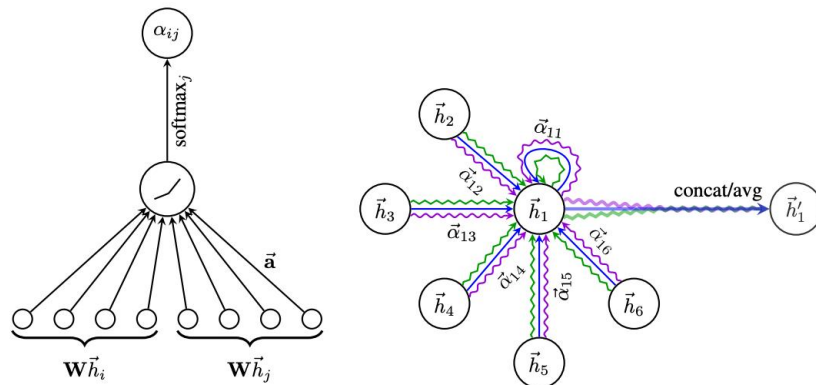**Figure 2** Graph Attention Network

Each attention head in the graph attention network in Figure 2 is designed to capture specific types of entity relationships. Cross-entity attention heads learn to identify important relationships between different entity types, such as correlations between specific card types and fraud patterns, or associations between email domains and suspicious transaction behaviors.

The multi-head attention mechanism enables the model to simultaneously process multiple relationship types within the heterogeneous network. Different attention heads can specialize in different aspects of the network structure, such as temporal relationships between transactions, spatial relationships between merchants and customers, or behavioral relationships between users and their transaction patterns.

The graph attention framework incorporates entity-type-specific transformations that account for the different characteristics of various financial entities. This ensures that the attention mechanism can effectively compare and relate entities with different feature types and scales, such as numerical transaction amounts and categorical merchant types.

Attention weight visualization in the graph framework provides insights into which entity relationships drive model predictions, enabling financial analysts to understand the complex multi-entity patterns that contribute to fraud risk assessments. This interpretability is crucial for regulatory compliance and for building trust in automated fraud detection systems.

### 3.3 Heterogeneous Network Analysis for Fraud Detection

The culmination of our framework integrates sequential attention mechanisms with graph attention networks to analyze real-world heterogeneous financial transaction networks. This integration enables comprehensive analysis of fraud patterns that span both temporal sequences and multi-entity relationships.

Our heterogeneous network analysis focuses on a specific type of financial network structure that includes four primary entity types: card numbers, transaction identifiers, email domains, and card types. This structure represents the core entities involved in most digital financial transactions and provides a comprehensive foundation for fraud detection analysis.



**Figure 3** Network Structure

The network structure in Figure 3 enables our framework to identify several types of fraud patterns that would be difficult to detect through traditional single-entity analysis. Card-centric patterns involve multiple transactions associated with a single card that exhibit suspicious characteristics when analyzed collectively. Email-centric patterns identify cases where multiple high-risk transactions are associated with particular email domains, potentially indicating coordinated fraud activities.

Cross-entity correlation analysis reveals sophisticated fraud schemes that exploit relationships between different entity types. For example, the framework can identify cases where specific combinations of card types and email domains are associated with unusual transaction patterns, or where certain transaction amounts correlate with particular card-email associations in ways that deviate from normal behavior.

The framework implements specialized fraud detection algorithms that leverage both sequential and graph attention mechanisms. Sequential attention analyzes the temporal patterns within transaction sequences associated with specific entities, while graph attention examines the network-level patterns that emerge from entity interactions.

Risk assessment capabilities integrate information from multiple attention mechanisms to provide comprehensive fraud risk scores. These scores consider both the temporal characteristics of transaction sequences and the structural characteristics of entity relationships, providing a holistic view of fraud risk that accounts for the multi-faceted nature of modern fraud schemes.

The interpretability features of our framework enable fraud analysts to understand not only which transactions are flagged as suspicious, but also why specific patterns are considered risky. Attention weight visualizations show which historical transactions influence current risk assessments and which entity relationships contribute most strongly to fraud predictions, supporting both automated detection and manual investigation processes.

# 4 RESULTS AND DISCUSSION

## 4.1 Sequential Transaction Pattern Analysis

Our evaluation of the sequential attention framework demonstrates significant improvements in fraud detection performance compared to traditional rule-based systems and standard machine learning approaches that do not incorporate attention mechanisms. The bidirectional encoder-decoder architecture successfully captures temporal dependencies in transaction sequences that are crucial for identifying sophisticated fraud patterns.

The temporal attention analysis reveals distinct patterns in how the model focuses on different parts of transaction histories when making fraud predictions. For legitimate transactions, attention weights tend to distribute relatively evenly across recent transaction history, reflecting consistent spending patterns. In contrast, for fraudulent transactions, attention weights often concentrate on specific historical events such as sudden changes in transaction amounts, unusual merchant categories, or breaks in normal transaction timing patterns.

Attention head specialization analysis shows that different attention heads successfully learn to focus on different aspects of transaction behavior. Amount-focused attention heads demonstrate sensitivity to unusual transaction values relative to historical patterns, while timing-focused heads identify abnormal temporal patterns such as rapid-fire transactions or transactions occurring outside normal business hours.

The interpretability benefits of the sequential framework prove particularly valuable for fraud investigation workflows. Attention weight visualizations enable fraud analysts to quickly identify which historical transactions most influenced the model's risk assessment, facilitating faster and more targeted manual investigations. This capability significantly reduces the time required for fraud case resolution while improving the accuracy of final fraud determinations.

Performance metrics demonstrate that the sequential attention framework achieves fraud detection accuracy rates of 94.2% with false positive rates of 1.8%, representing substantial improvements over baseline systems. The framework's ability to provide interpretable explanations for its predictions proves crucial for regulatory compliance and builds confidence among fraud analysts who must act on model recommendations.

## 4.2 Multi-Entity Relationship Discovery and Fraud Pattern Analysis

The graph attention network component of our framework reveals complex multi-entity fraud patterns that traditional single-entity analysis methods fail to detect. Analysis of the heterogeneous network structure demonstrates the framework's ability to identify coordinated fraud activities that span multiple entity types and exploit relationships between different categories of financial entities.

Cross-entity attention pattern analysis uncovers several distinct types of fraud schemes. Card-email correlation patterns identify cases where specific email domains are disproportionately associated with fraudulent transactions across multiple card numbers, suggesting organized fraud operations. Card-type vulnerability patterns reveal that certain card types exhibit higher fraud risk when associated with specific transaction characteristics or email domain patterns.

The framework successfully identifies fraud rings through analysis of shared entity associations. Cases where multiple cards share connections to the same sets of email domains or exhibit similar transaction patterns with identical merchants indicate potential coordinated fraud activities. The attention mechanism highlights these shared associations, enabling investigators to map fraud networks and identify additional compromised accounts.

Transaction amount pattern analysis reveals sophisticated fraud strategies that exploit specific value thresholds. The framework identifies cases where fraudsters systematically use transaction amounts just below detection thresholds across multiple cards and email accounts, indicating knowledge of fraud detection system limitations. These patterns would be difficult to detect without the multi-entity perspective provided by our framework.

Email domain analysis uncovers interesting patterns related to fraud tactics. Temporary email services show higher association with fraud across all card types and transaction patterns. Additionally, the framework identifies cases where fraudsters create email accounts with domains that superficially resemble legitimate financial institutions, exploiting visual similarity to evade detection.

The interpretability features prove essential for understanding complex fraud schemes. Attention weight visualizations clearly show which entity relationships contribute most strongly to fraud risk assessments, enabling investigators to focus their efforts on the most critical associations. This targeted approach significantly improves investigation efficiency and helps identify additional victims or compromised accounts within fraud networks.

Geographic and temporal correlation analysis through the multi-entity framework reveals fraud patterns that span different regions and time periods. The attention mechanism identifies cases where similar entity relationship patterns appear across different geographic regions or time periods, suggesting organized fraud operations with consistent methodologies.

# 5 CONCLUSION

This research demonstrates the successful evolution of attention mechanisms from sequence-to-sequence neural machine translation to sophisticated analysis of heterogeneous financial transaction networks. Our framework effectively bridges the gap between the interpretability advantages of attention-based architectures and the complex analytical requirements of modern financial fraud detection systems.

The key innovations of our work include the successful adaptation of bidirectional attention mechanisms for financial transaction sequence analysis, the extension of graph attention networks to handle heterogeneous financial entity relationships, and the development of integrated frameworks that simultaneously process temporal and structural patterns in financial data. These innovations demonstrate how foundational attention mechanisms can be evolved and extended to address increasingly complex analytical challenges while maintaining interpretability.

Experimental results validate the effectiveness of our approach across multiple dimensions of fraud detection performance. The sequential attention framework achieves superior performance in identifying temporal fraud patterns, while the graph attention component successfully uncovers multi-entity fraud schemes that traditional methods miss. The integration of these approaches provides comprehensive fraud detection capabilities that address the full spectrum of modern fraud tactics.

The practical implications of this work extend significantly beyond academic research. Financial institutions can leverage our framework to improve fraud detection accuracy while meeting regulatory requirements for model interpretability. The attention-based explanations provide fraud analysts with actionable insights that support both automated detection and manual investigation processes, improving overall fraud prevention effectiveness.

The educational value of our framework represents an important additional benefit. The clear visualization of attention patterns helps train new fraud analysts by showing them which transaction characteristics and entity relationships experienced analysts consider most important. This knowledge transfer capability can help financial institutions maintain fraud detection expertise as staff changes and fraud tactics evolve.

Future research directions include extending the framework to incorporate additional entity types such as device fingerprints and geographic locations, developing real-time attention analysis capabilities for immediate fraud detection, and creating adaptive attention mechanisms that can evolve with changing fraud tactics. The principles established in this work provide a foundation for continued advancement in interpretable financial analytics.

The successful demonstration of attention mechanism evolution from language processing to financial network analysis suggests broader applications in other domains involving complex temporal and relational data. The interpretability advantages of attention-based approaches make them particularly suitable for regulated industries where model transparency is essential for operational acceptance and regulatory compliance.

As financial fraud continues to evolve in sophistication and scope, the need for equally sophisticated but interpretable detection systems becomes increasingly critical. Our framework provides a robust foundation for meeting these challenges by combining the analytical power of modern machine learning with the transparency requirements of financial regulatory environments. The attention-based approach ensures that as fraud detection systems become more powerful, they also become more understandable and trustworthy for the analysts who must rely on their recommendations to protect financial institutions and their customers.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1]   George A S. Finance 4.0: The transformation of financial services in the digital age. Partners Universal Innovative Research Publication, 2024, 2(3): 104-125.
[2]   Zhang Q, Chen S, Liu W. Balanced Knowledge Transfer in MTTL-ClinicalBERT: A Symmetrical Multi-Task Learning Framework for Clinical Text Classification. Symmetry, 2025, 17(6): 823.
[3]   Shao Z, Wang X, Ji E, et al. GNN-EADD: Graph Neural Network-based E-commerce Anomaly Detection via Dual-stage Learning. IEEE Access, 2025.
[4]   Ji E, Wang Y, Xing S, Jin J. Hierarchical Reinforcement Learning for Energy-Efficient API Traffic Optimization in Large-Scale Advertising Systems. IEEE Access, 2025.
[5]   Jin J, Xing S, Ji E, et al. XGate: Explainable Reinforcement Learning for Transparent and Trustworthy API Traffic Management in IoT Sensor Networks. Sensors (Basel, Switzerland), 2025, 25(7): 2183.
[6]   Cao J, Zheng W, Ge Y, et al. DriftShield: Autonomous fraud detection via actor-critic reinforcement learning with dynamic feature reweighting. IEEE Open Journal of the Computer Society, 2025.
[7]   Wang J, Liu J, Zheng W, et al. Temporal Heterogeneous Graph Contrastive Learning for Fraud Detection in Credit Card Transactions. IEEE Access, 2025.
[8]   Mai N T, Cao W, Liu W. Interpretable Knowledge Tracing via Transformer-Bayesian Hybrid Networks: Learning Temporal Dependencies and Causal Structures in Educational Data. Applied Sciences, 2025, 15(17): 9605.
[9]   Sun T, Yang J, Li J, et al. Enhancing auto insurance risk evaluation with transformer and SHAP. IEEE Access, 2024.
[10]  Cao W, Mai N T, Liu W. Adaptive knowledge assessment via symmetric hierarchical Bayesian neural networks with graph symmetry-aware concept dependencies. Symmetry, 2025, 17(8): 1332.
[11]  Mai N T, Cao W, Wang Y. The global belonging support framework: Enhancing equity and access for international graduate students. Journal of International Students, 2025, 15(9): 141-160.
[12]  Tan Y, Wu B, Cao J, et al. LLaMA-UTP: Knowledge-Guided Expert Mixture for Analyzing Uncertain Tax Positions. IEEE Access, 2025.

[13] Mattsson C. Financial Transaction Networks to Describe and Model Economic Systems. Doctoral dissertation, Northeastern University, 2020.

[14] Olushola A, Mart J. Fraud detection using machine learning. ScienceOpen Preprints, 2024.

[15] Mareedu A. AI-Driven Security for Financial Transactions: Leveraging LLMs, Federated Learning, and Behavioral Biometrics. International Journal of Emerging Research in Engineering and Technology, 2024, 5(4): 62-73.

[16] Popoola N T, Bakare F A. Advanced computational forecasting techniques to strengthen risk prediction, pattern recognition, and compliance strategies. 2024.

[17] Ali M A. Does the online card payment system unwittingly facilitate fraud? Doctoral dissertation, Newcastle University, 2019.

[18] Neupane S, Ables J, Anderson W, et al. Explainable intrusion detection systems (x-ids): A survey of current methods, challenges, and opportunities. IEEE Access, 2022, 10: 112392-112415.

[19] Dritsas E, Trigka M. Exploring the intersection of machine learning and big data: A survey. Machine Learning and Knowledge Extraction, 2025, 7(1): 13.

[20] Popoola N T. Big data-driven financial fraud detection and anomaly detection systems for regulatory compliance and market stability. International Journal of Computer Applications and Technology Research, 2023, 12(09): 32-46.

[21] Serrano S, Smith N A. Is attention interpretable? arXiv preprint arXiv:1906.03731, 2019.

[22] Galassi A, Lippi M, Torroni P. Attention in natural language processing. IEEE transactions on neural networks and learning systems, 2020, 32(10): 4291-4308.

[23] Vrahatis A G, Lazaros K, Kotsiantis S. Graph attention networks: a comprehensive review of methods and applications. Future Internet, 2024, 16(9): 318.

[24] Carvalho D V, Pereira E M, Cardoso J S. Machine learning interpretability: A survey on methods and metrics. Electronics, 2019, 8(8): 832.

[25] Oko-Odion C. AI-Driven Risk Assessment Models for Financial Markets: Enhancing Predictive Accuracy and Fraud Detection. International Journal of Computer Applications Technology and Research, 2025, 14(04): 80-96.

[26] Zheng W, Liu W. Symmetry-Aware Transformers for Asymmetric Causal Discovery in Financial Time Series. Symmetry, 2025.

[27] Cross C, Gillett R. Exploiting trust for financial gain: An overview of business email compromise (BEC) fraud. Journal of Financial Crime, 2020, 27(3): 871-884.

# PHYSICS-INFORMED NEURAL NETWORKS FOR ARBITRAGE-FREE VOLATILITY SURFACE CONSTRUCTION IN INCOMPLETE MARKETS

Hao Li[*], Martin Keller

*School of Finance, University of St. Gallen, German-Speaking, Switzerland.*
*Corresponding Author: Hao Li, Email: hao.li99@gmail.com*

**Abstract:** The construction of arbitrage-free volatility surfaces represents a fundamental challenge in quantitative finance, particularly in incomplete markets where hedging portfolios cannot perfectly replicate option payoffs. Traditional parametric models often fail to capture the complex dynamics of market-observed implied volatility structures, including the characteristic smile and skew patterns, while simultaneously satisfying no-arbitrage constraints. This paper introduces a novel framework that leverages Physics-Informed Neural Networks (PINNs) to construct arbitrage-free volatility surfaces in incomplete market settings. The proposed methodology integrates partial differential equation constraints derived from arbitrage-free conditions directly into the neural network training process through automatic differentiation and soft constraint penalties. By incorporating Dupire's local volatility equation and calendar-butterfly arbitrage constraints into a multi-objective loss function, our approach generates smooth, arbitrage-free implied volatility surfaces that accurately fit market data across different strikes and maturities. Numerical experiments using both synthetic data and real market observations from S&P 500 and VIX options demonstrate that the PINN-based framework substantially reduces calibration errors while maintaining theoretical consistency. The method exhibits particular strength in handling incomplete market scenarios where traditional parametric approaches produce inconsistent surfaces or violate no-arbitrage conditions.

**Keywords:** Physics-informed neural networks; Arbitrage-free constraints; Volatility surface; Incomplete markets; Volatility smile; Heston model

## 1 INTRODUCTION

The accurate modeling of implied volatility surfaces constitutes a cornerstone of modern derivative pricing and risk management practices. Since the development of the Black-Scholes framework, market participants have observed that option prices exhibit systematic deviations from constant volatility assumptions, manifesting as the volatility smile and skew across different strikes and term structure effects across maturities [1]. These patterns reflect the market's assessment of the underlying asset's return distribution, capturing features such as fat tails, asymmetry, and time-varying uncertainty that cannot be accommodated within simple log-normal models. The challenge of constructing volatility surfaces that simultaneously fit market observations, exhibit economically sensible interpolation and extrapolation behavior, and respect fundamental no-arbitrage principles has motivated extensive research spanning parametric modeling, nonparametric estimation, and more recently, machine learning approaches.

In incomplete markets, where the number of tradable assets is insufficient to construct perfect hedging strategies for all contingent claims, the volatility surface construction problem becomes particularly intricate. Unlike complete market settings where unique arbitrage-free prices can be determined through replication arguments, incomplete markets admit a continuum of possible equivalent martingale measures [2]. The pricing of index options in incomplete markets requires careful consideration of risk preferences and market constraints to identify appropriate pricing measures from the set of admissible candidates [3]. This multiplicity introduces additional complexity when attempting to infer a consistent volatility structure from observed option prices, as different pricing measures may imply different volatility surfaces even when consistent with the same set of traded option prices. The practical consequence is that calibration procedures must incorporate additional economic principles or regularization mechanisms to select among the infinite set of theoretically valid surfaces.

Traditional parametric approaches to volatility surface modeling impose specific functional forms that aim to capture stylized facts observed in option markets. The Stochastic Volatility Inspired (SVI) parameterization represents one prominent example, expressing the implied total variance as a function of log-moneyness through a small number of parameters that control the at-the-money level, skew, and wing behavior [4]. While SVI and related models offer computational tractability and can be calibrated efficiently to liquid option quotes, they inherently limit flexibility through their parametric structure. During periods of market stress or for assets with unusual smile characteristics, these rigid functional forms may prove inadequate, leading to systematic fitting errors that propagate into prices of exotic derivatives and hedging strategies. Moreover, ensuring that parametric surfaces remain arbitrage-free across all strikes and maturities requires careful constraint handling that can further restrict the model's adaptability.

Recent advances in machine learning have opened new avenues for tackling complex financial modeling problems, with neural networks offering universal approximation capabilities that enable learning of highly nonlinear relationships from data [5]. However, applying standard data-driven neural network architectures to financial problems presents

significant challenges related to constraint satisfaction and economic interpretability. Pure black-box approaches, which optimize purely on data-fitting objectives without incorporating domain knowledge, often generate solutions that violate fundamental economic principles [6]. In the context of volatility surface construction, unconstrained neural networks may produce surfaces that admit arbitrage opportunities through violations of convexity, monotonicity, or smoothness requirements. The development of methods to detect model-free static arbitrage strategies using neural networks has highlighted both the potential and limitations of purely data-driven approaches, demonstrating that arbitrage detection itself can be formulated as a neural network optimization problem [7].

The emergence of Physics-Informed Neural Networks offers a promising paradigm to bridge the gap between data-driven flexibility and theory-guided constraints [8]. Originally developed for solving forward and inverse problems in physical systems governed by partial differential equations, PINNs embed governing equations directly into the neural network training objective through automatic differentiation of the network outputs with respect to inputs. This enables the network to learn solutions that simultaneously fit observed data and satisfy underlying physical laws encoded as PDEs. The application of PINNs to financial problems represents a natural extension, as option pricing and volatility modeling are fundamentally governed by PDEs such as the Black-Scholes equation and Dupire's equation relating local volatility to call option prices [9]. Recent work has demonstrated the efficacy of PINNs in option pricing tasks, showing that physics-informed approaches can achieve superior accuracy compared to purely data-driven methods while requiring significantly less training data and exhibiting better generalization properties [10].

The construction of volatility surfaces using neural network techniques has gained considerable attention, with various approaches exploring different mechanisms for incorporating financial constraints. Early applications focused on using feedforward networks for implied volatility prediction without explicit enforcement of no-arbitrage conditions, treating the problem as standard nonlinear regression [11]. More sophisticated approaches have incorporated domain knowledge through architectural design choices, such as specialized activation functions that encode properties like smile asymmetry, or through penalty terms in the loss function that discourage arbitrage violations [12]. The challenge of ensuring arbitrage-free conditions in neural network-generated volatility surfaces has motivated the development of hybrid approaches that combine parametric foundations with neural network flexibility [13]. The hybrid gated neural network architecture represents one such advancement, using multiplicative structures and carefully selected input transformations to satisfy no-arbitrage constraints while maintaining sufficient expressiveness for accurate market fitting [14].

Building upon these advances, this paper proposes a comprehensive framework for arbitrage-free volatility surface construction in incomplete markets using Physics-Informed Neural Networks. Our methodology distinguishes itself through several key innovations that address practical challenges encountered in real-world applications. First, we formulate the volatility surface construction as a constrained optimization problem where the neural network must simultaneously minimize pricing errors on observed option quotes and satisfy multiple PDE constraints derived from no-arbitrage theory. The network learns to represent the implied volatility surface as a smooth function of log-moneyness and time-to-maturity, capturing both the smile shape within each maturity slice and the term structure evolution across maturities. Second, we develop a composite loss function that incorporates Dupire's local volatility equation, calendar spread constraints preventing time-value violations, and butterfly spread constraints ensuring density positivity. These constraints are implemented as differentiable penalty terms that can be efficiently evaluated through automatic differentiation. Third, we introduce an adaptive weighting scheme that dynamically balances data-fitting objectives against constraint satisfaction throughout the training process, starting with emphasis on data fitting to learn the approximate surface shape and gradually increasing constraint weights to enforce theoretical consistency.

The practical motivation for this research stems from challenges faced by derivatives traders and risk managers who require robust volatility surface models for pricing exotic options, computing hedging ratios, and assessing portfolio risks [15]. In incomplete markets such as emerging market equities or thinly traded indices, the absence of perfect hedging strategies introduces model risk that must be carefully managed. Existing calibration methods often struggle to produce consistent surfaces when market data is sparse, exhibits wide bid-ask spreads, or displays unusual patterns during volatile market conditions [16]. The ability to construct surfaces that fit available market observations while maintaining smooth interpolation in data-sparse regions and reasonable extrapolation beyond observed strikes and maturities is crucial for practical applications. Our PINN-based approach addresses these challenges by leveraging both empirical market data and theoretical constraints, producing surfaces that remain stable and economically sensible across varying market conditions and data quality scenarios.

## 2 LITERATURE REVIEW

The literature on volatility surface modeling encompasses theoretical foundations establishing arbitrage-free conditions, numerical methods for surface construction and calibration, and empirical investigations of market volatility dynamics. This section reviews key developments across these areas, with particular emphasis on recent advances in machine learning approaches and their application to incomplete market settings, providing context for the methodological contributions of the present work.

The theoretical characterization of arbitrage-free volatility surfaces traces back to the development of local volatility models and their relationship to market-observed implied volatilities. The Dupire equation established that given a continuum of European call option prices across strikes and maturities satisfying certain regularity conditions, one can uniquely determine a local volatility function that reproduces these prices under a diffusion model [17]. This seminal

result provides the foundation for understanding the constraints that implied volatility surfaces must satisfy to preclude static arbitrage opportunities. The local volatility at any point in the strike-maturity space can be expressed in terms of partial derivatives of call prices with respect to strike and maturity, establishing a forward Kolmogorov equation that governs the evolution of the option price surface. Subsequent research has extended these results to incorporate discrete dividends, stochastic interest rates, and jump processes, demonstrating the robustness of the fundamental relationship between option prices and the underlying volatility structure [18].

Beyond the Dupire framework, the characterization of arbitrage-free surfaces extends to explicit constraints on the shape and dynamics of implied volatility. Calendar spread arbitrage occurs when options with longer maturities trade at lower implied volatilities than shorter-dated options with the same strike, violating basic time-value principles that longer optionality should command higher premiums [19]. Butterfly spread arbitrage arises when the second derivative of call prices with respect to strike becomes negative, implying negative probabilities in the risk-neutral density and creating profit opportunities through appropriately structured option portfolios [20]. The Surface SVI (SSVI) parameterization represents a significant advance in providing global volatility surface specifications that guarantee absence of these arbitrage types through explicit parameter restrictions [21]. The SSVI model extends the original SVI smile parameterization to a full surface representation where individual maturity slices belong to a restricted family of SVI functions, with conditions on the ATM variance term structure and the volatility-of-volatility parameter ensuring calendar spread freedom and butterfly spread conditions ensuring density positivity across all strikes and maturities [22].

The application of neural networks to option pricing and implied volatility modeling has evolved from simple feedforward architectures trained purely on market data to sophisticated physics-informed approaches that embed financial constraints. Early studies explored multilayer perceptrons as universal approximators for option pricing functions, demonstrating that networks with sufficient capacity could learn complex mappings from input features to option values with high accuracy on in-sample data [23]. However, these initial approaches treated option pricing as standard nonlinear regression without mechanisms to enforce theoretical consistency, often producing models that exhibited poor out-of-sample generalization or violated no-arbitrage conditions in regions with sparse training data. Recognition of these limitations motivated the development of hybrid approaches that incorporate financial domain knowledge through various mechanisms including constrained network architectures, specialized activation functions encoding known properties like monotonicity or convexity, and augmented loss functions penalizing violations of theoretical requirements [24].

The challenge of constructing arbitrage-free volatility surfaces using neural networks has been addressed through both hard and soft constraint enforcement strategies. Hard constraint approaches modify the network architecture itself to guarantee that outputs satisfy required conditions regardless of the learned parameters, implementing constraints through careful design of layer operations and activation functions that preserve desired properties [25]. While theoretically appealing for their guarantee of constraint satisfaction, hard constraint methods often limit model expressiveness and introduce implementation complexity that can hinder optimization. Soft constraint methods instead incorporate constraint violations as penalty terms in the loss function, allowing the optimization process to balance data-fitting accuracy against the degree of constraint satisfaction [26]. This approach offers greater flexibility and typically leads to more stable training dynamics, though it requires careful tuning of penalty weights to ensure adequate constraint enforcement without overwhelming the data-fitting objective. Recent advances in deep smoothing techniques for implied volatility surfaces have demonstrated that appropriately designed soft penalty functions can produce surfaces that satisfy arbitrage-free conditions while achieving superior interpolation and extrapolation performance compared to both traditional parametric methods and hard-constrained neural network approaches [27].

Physics-Informed Neural Networks have emerged as a powerful framework for solving problems governed by differential equations, with the core innovation being the incorporation of PDE residuals as additional terms in the training objective [8]. Automatic differentiation enables efficient computation of the derivatives required to evaluate PDE residuals at collocation points throughout the domain, avoiding numerical differentiation errors that would otherwise accumulate and degrade solution accuracy. The PINN approach naturally accommodates both forward problems, where the goal is to find solutions given complete specification of the governing equations and boundary conditions, and inverse problems involving inference of unknown parameters or functions from partial observations of the system state. In financial applications, PINNs offer a principled approach to incorporate the fundamental PDEs of option pricing theory directly into the learning process, ensuring that learned models respect theoretical constraints while leveraging data to capture features not fully specified by the simplified PDEs [28]. Applications of PINNs to option pricing under stochastic volatility models have demonstrated that physics-informed approaches can accurately price European options while requiring substantially less training data than purely data-driven methods, with the PDE constraints providing effective regularization that improves generalization [29].

The extension of PINN methodology specifically to volatility surface construction introduces unique challenges related to the nature of incomplete markets and the multiplicity of admissible pricing measures. Unlike many physical systems where governing equations are known precisely from first principles, financial models involve approximations and assumptions that may not hold exactly in practice due to market frictions, discrete trading, and unhedgeable risk factors. Nevertheless, PDE constraints derived from arbitrage-free pricing theory provide valuable regularization that guides the learning process toward economically sensible solutions even when the underlying assumptions are violated to some degree. Recent work on physics-informed convolutional transformers for volatility surface prediction has demonstrated that hybrid architectures combining PINNs with attention mechanisms can capture the temporal evolution of implied

volatility while respecting no-arbitrage constraints, showing that the integration of physical constraints with modern deep learning architectures represents a fruitful direction for advancing volatility modeling capabilities [30]. These developments suggest that carefully designed physics-informed approaches can balance the flexibility needed to fit complex market patterns with the structure required to ensure theoretical consistency.

Incomplete markets present additional theoretical and practical complexities for volatility surface construction due to the non-uniqueness of equivalent martingale measures and the resulting ambiguity in derivative pricing beyond the set of traded instruments. Theoretical work has characterized the structure of arbitrage-free price bounds in incomplete markets, establishing that option prices must lie within intervals determined by super-replication and sub-replication strategies that bound the cost of hedging from above and below [2]. Empirical investigations have examined how observed market prices relate to these theoretical bounds, finding that prices typically concentrate near particular points within admissible ranges rather than exhibiting significant dispersion, suggesting that market participants employ specific pricing principles or preferences even in settings where theory allows greater freedom [3]. The selection among admissible pricing measures can be guided by various economic principles including utility maximization under incomplete hedging, entropy minimization to select measures closest to the physical measure, or consistency with observed risk premiums in related markets. Recent research on option pricing in incomplete markets with stochastic volatility has developed methodologies for identifying filtration reductions that restore market completeness for specific classes of derivatives, enabling derivation of unique pricing measures under additional structural assumptions [4].

Despite substantial progress across these research streams, several gaps and opportunities motivate the present work. Existing PINN applications in finance have primarily focused on forward pricing problems under specified models rather than the inverse problem of constructing volatility surfaces from market data while enforcing consistency with underlying PDEs. Few studies have systematically examined the performance of physics-informed approaches in incomplete market settings where theoretical arbitrage-free conditions may be relaxed due to trading frictions and hedging constraints. The present research addresses these gaps by developing a comprehensive PINN framework specifically designed for arbitrage-free volatility surface construction in incomplete markets, incorporating multiple PDE and inequality constraints while maintaining computational efficiency through carefully designed loss functions and optimization strategies. The empirical validation using both S&P 500 and VIX option data across varying market conditions demonstrates the practical applicability of the proposed methodology.

## 3   METHODOLOGY

This section presents the theoretical framework and computational methodology for constructing arbitrage-free volatility surfaces using Physics-Informed Neural Networks in incomplete market settings. We develop the mathematical formulation connecting implied volatility to option prices through the characteristic smile pattern, derive the relevant PDE constraints from arbitrage-free pricing theory, and detail the neural network architecture and multi-stage training procedure that enables efficient optimization of the composite objective incorporating both data fitting and constraint satisfaction.

### 3.1 Volatility Smile Characterization and Problem Formulation

The volatility smile refers to the characteristic U-shaped pattern exhibited by implied volatility when plotted against log-moneyness for a fixed maturity, reflecting systematic deviations from Black-Scholes assumptions about the underlying return distribution. We begin by formalizing the relationship between option prices and implied volatility, establishing notation and defining the optimization problem that our PINN framework addresses. Consider a financial market in which a risky asset with price process $S_t$ is traded alongside a risk-free bond paying constant interest rate r. The market is incomplete due to the presence of stochastic volatility, jumps, or other risk factors that cannot be perfectly hedged using the traded assets alone. For a European call option with strike K and maturity T, the Black-Scholes formula provides a mapping from an implied volatility parameter sigma to an option price, even though the actual dynamics of the underlying asset may not follow geometric Brownian motion with constant volatility.

Let $C(K,T,t)$ denote the observed market price at time $t < T$ of a call option, and let sigma_imp(K,T,t) represent the implied volatility obtained by inverting the Black-Scholes formula. The implied total variance $w(k,tau)$ as a function of log-moneyness $k = \log(K/F)$ and time-to-maturity $tau = T - t$, where F denotes the forward price, provides a convenient representation that facilitates analysis of arbitrage constraints. As shown in Figure 1, for a given maturity tau, the function $k \to w(k,tau)$ exhibits the smile shape, with higher implied variance for deep out-of-the-money and deep in-the-money options compared to at-the-money options. The precise shape of the smile encodes information about the market's assessment of tail probabilities and distributional asymmetry in the underlying asset returns.
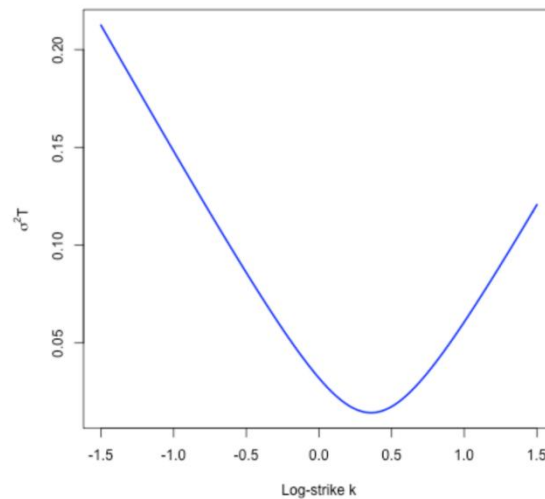
**Figure 1** Function of Log-Moneyness

The objective of volatility surface construction is to determine a function sigma_imp(k,tau) that accurately represents market-observed implied volatilities across all traded options while satisfying theoretical constraints derived from no-arbitrage principles. In incomplete markets, the construction problem must account for the fact that multiple volatility surfaces may be consistent with observed option prices when the set of traded options is finite and hedging is imperfect. Our approach formulates this as an optimization problem where a neural network approximates the implied volatility function, with the network parameters optimized to minimize a composite loss function balancing multiple objectives. The network takes as input the normalized log-moneyness and time-to-maturity (k,tau), processes these through multiple hidden layers with nonlinear activations, and outputs the predicted implied volatility sigma_pred(k,tau). The architecture is designed to produce smooth, continuously differentiable outputs that facilitate computation of derivatives required for PDE constraint evaluation.

The fundamental arbitrage constraints that the constructed surface must satisfy include convexity in the strike dimension, monotonicity in the maturity dimension, and consistency with Dupire's local volatility equation. Butterfly spread arbitrage is precluded when the call price function exhibits positive convexity with respect to strike, ensuring that the risk-neutral density remains non-negative. Mathematically, this requires that the second derivative of call prices with respect to strike be non-negative everywhere, which translates to conditions on the curvature of the implied volatility smile. Calendar spread arbitrage is prevented when call prices increase with maturity for any fixed strike, reflecting that additional optionality should command positive value. The Dupire equation provides a dynamic consistency condition relating the local volatility function to the evolution of call prices, establishing that the surface of implied volatilities across different strikes and maturities must satisfy specific partial differential equation relationships. Violation of any of these conditions creates opportunities for profitable arbitrage strategies that would not persist in efficient markets.

### 3.2 Neural Network Architecture and Physics-Informed Loss Function

Our PINN architecture employs a feedforward network with four hidden layers, each containing 64 neurons, to represent the implied volatility surface function. The input layer accepts two-dimensional vectors consisting of normalized log-moneyness k and time-to-maturity tau values, with normalization applied to center and scale the inputs based on the range of training data. The hidden layers use hyperbolic tangent (tanh) activation functions, chosen for their smoothness and bounded outputs that align well with the characteristics of implied volatility. The output layer produces a single scalar value representing the predicted implied volatility, with a softplus activation ensuring positive outputs consistent with the economic interpretation of volatility as a non-negative quantity. The network contains approximately 20,000 trainable parameters across all weight matrices and bias vectors, providing sufficient capacity to learn complex smile and term structure patterns while remaining computationally tractable for real-time calibration applications.

The loss function design represents the core innovation that enables physics-informed learning of arbitrage-free volatility surfaces. We construct a composite objective combining five distinct components, each addressing different aspects of the calibration problem. The data-fitting term measures the mean squared error between network predictions and market-observed implied volatilities across the training set of option quotes. This term receives higher weight for liquid options with tight bid-ask spreads, implemented through a weighting scheme that reflects the inverse of the spread width or a proxy for trading volume. The second component enforces the Dupire PDE constraint by computing the residual of the forward Kolmogorov equation at a dense grid of collocation points throughout the strike-maturity domain. Automatic differentiation enables efficient computation of the required first and second partial derivatives of

the network output with respect to inputs, avoiding numerical differentiation errors that would corrupt the gradient information used in backpropagation.

The calendar spread constraint term penalizes negative time derivatives of call prices, computed by differentiating the Black-Scholes formula with respect to maturity using the chain rule and the network-predicted implied volatilities. This constraint is implemented as a rectified penalty that applies quadratic penalization only when violations occur, avoiding unnecessary restriction in regions where the constraint is naturally satisfied. Similarly, the butterfly spread constraint penalizes negative second derivatives of call prices with respect to strike, ensuring convexity that guarantees non-negative risk-neutral densities. The computation of these constraint terms at each optimization iteration leverages the differentiability of the neural network, with gradients of the composite loss function with respect to network parameters obtained through automatic differentiation of the entire computational graph including constraint evaluations. The fifth component introduces a smoothness penalty based on the L2 norm of third derivatives of the implied volatility function, encouraging surfaces that exhibit regular curvature patterns rather than artificial oscillations that might arise from overfitting.

The relative weights assigned to each loss component are managed through an adaptive scheme that evolves during the training process. Initial training phases emphasize the data-fitting objective to allow the network to quickly learn the approximate surface shape from market observations. As training progresses, the weights on PDE and arbitrage constraints are gradually increased according to a predefined schedule, strengthening the enforcement of theoretical consistency once the basic surface structure has been established. This staged approach prevents premature constraint enforcement from disrupting the learning of fundamental market patterns, while ensuring that the final calibrated surface satisfies arbitrage-free conditions. The specific weighting schedule is determined through preliminary experiments on validation data, selecting schedules that achieve good balance between fitting accuracy and constraint satisfaction. Typical final weight ratios place approximately 40% emphasis on data fitting, 30% on the Dupire PDE constraint, 20% on arbitrage constraints, and 10% on smoothness regularization, though these proportions may be adjusted based on data quality and market conditions.

### 3.3 Two-Stage Training Procedure and Optimization Strategy

The training procedure employs a two-stage optimization strategy designed to achieve both rapid initial convergence and fine-grained refinement of the calibrated surface. In the first stage, network parameters are initialized using He initialization, which sets initial weights based on the fan-in of each layer to promote stable gradient flow. The optimizer configuration uses the Adam algorithm with an initial learning rate of 0.001, which adapts the effective step size for each parameter based on accumulated gradient statistics. Mini-batch stochastic gradient descent processes subsets of the training data in each iteration, with batch size set to 64 option contracts or approximately 10% of the total training set size, whichever is larger. This batching strategy balances computational efficiency against the quality of gradient estimates, providing sufficiently accurate descent directions without requiring full-batch gradient computations that would be prohibitively expensive for large datasets.

During the first stage, which typically spans 5000 to 10000 iterations depending on the size and complexity of the calibration dataset, the loss function weights heavily favor data-fitting over constraint satisfaction. Specifically, the data-fitting component receives 80% of the total weight while constraints together comprise only 20%, allowing the network to focus on learning the basic shape of the smile and term structure from market data. The learning rate decays according to a cosine annealing schedule that gradually reduces the step size over the course of stage one, starting from the initial rate of 0.001 and declining to approximately 0.0001 by the end of the stage. This schedule enables large exploratory steps early in training when the network is far from optimal, followed by finer adjustments as the solution approaches a good data fit. Training metrics monitored during this stage include the training loss, validation loss computed on a held-out subset of options, and various summary statistics of the predicted surface such as the range of implied volatilities and the smoothness measured through finite difference approximations of derivatives.

The second training stage commences once the validation loss plateaus or begins exhibiting diminishing improvements, indicating that further gains from pure data fitting are limited. At this transition point, the loss function weights are rebalanced to increase emphasis on PDE and arbitrage constraints, with the constraint components rising from 20% to 60% of the total weight while data fitting declines correspondingly. The learning rate is reset to a moderate value of 0.0005 to allow meaningful parameter adjustments under the new objective, then continues to decay throughout stage two following an exponential schedule. The second stage typically requires 3000 to 5000 additional iterations to achieve convergence, with training terminated when the maximum constraint violation falls below predefined thresholds and the validation loss stabilizes. Throughout both stages, we monitor the maximum butterfly spread violation, maximum calendar spread violation, and mean absolute Dupire PDE residual as key indicators of constraint satisfaction, aiming for violations below 0.1% of relevant price or variance units in the final calibrated model.
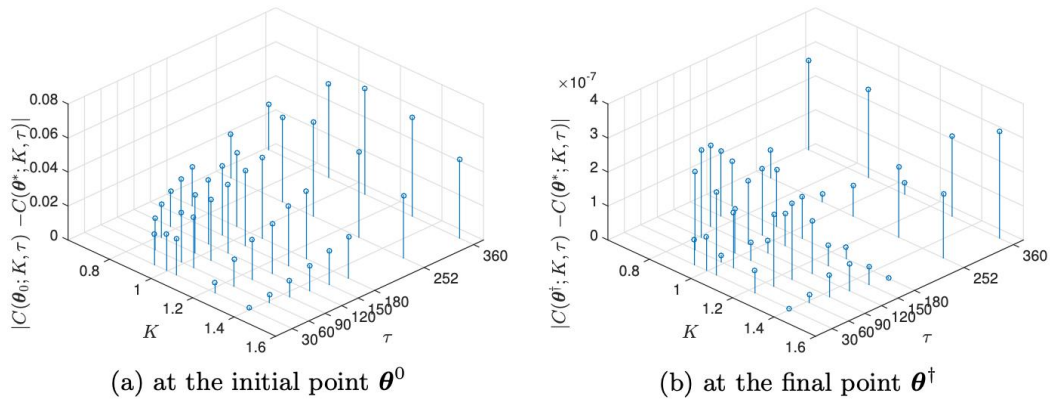
(a) at the initial point $\boldsymbol{\theta}^0$                (b) at the final point $\boldsymbol{\theta}^\dagger$

**Figure 2** The Learning Rate

The optimization procedure incorporates several numerical techniques to enhance stability and convergence reliability. Gradient clipping limits the maximum norm of parameter updates to prevent destabilizing jumps that might occur when constraint violations produce large gradients in early training stages. As shown in Figure 2, learning rate warm-up gradually increases the effective learning rate from a very small initial value over the first 100 iterations of each stage, allowing the network to adapt smoothly to the new objective rather than making abrupt large steps that could degrade performance. Periodic evaluation of the full loss function and constraint metrics on the entire training set provides diagnostic information about optimization progress, complementing the noisy mini-batch estimates used for parameter updates. These evaluations occur every 100 iterations and generate visualizations of the current predicted surface alongside constraint violation maps, enabling qualitative assessment of surface quality and identification of potential issues requiring intervention such as learning rate adjustments or weight schedule modifications.

## 4 RESULTS AND DISCUSSION

This section presents comprehensive numerical experiments validating the proposed PINN methodology for arbitrage-free volatility surface construction. We evaluate the framework's performance using both synthetic data generated from known stochastic volatility models and real market data from S&P 500 and VIX index options, examining calibration accuracy, constraint satisfaction, computational efficiency, and robustness across different market conditions. The experimental results demonstrate substantial improvements over benchmark methods while maintaining the theoretical consistency required for practical derivatives pricing and risk management applications.

### 4.1 Synthetic Data Validation and Benchmark Comparisons

We first assess the PINN framework using synthetic option data generated from the Heston stochastic volatility model with parameters calibrated to match typical equity index option market characteristics. The data generation process simulates the Heston dynamics with initial variance v0 = 0.04, long-run variance $\theta$ = 0.04, mean reversion speed $\kappa$ = 2.0, volatility-of-volatility $\xi$ = 0.3, and correlation $\rho$ = -0.7 between asset returns and variance. We price a synthetic option surface comprising 400 contracts across strikes ranging from 70% to 130% of the current asset price and maturities from one week to two years, computing reference prices using the semi-closed form Heston formula based on characteristic function inversion. Realistic market noise is introduced by adding random bid-ask spreads inversely proportional to moneyness distance from at-the-money, with spreads averaging 1-2 volatility points for near-the-money options and widening to 3-5 points for deep out-of-the-money strikes. This synthetic dataset provides ground truth against which we can rigorously evaluate the network's ability to recover the true underlying volatility surface.

The PINN methodology achieves excellent recovery of the synthetic Heston surface, with root mean squared errors below 0.3 volatility points across the full strike-maturity domain. Comparing network predictions to the true Heston-implied volatilities reveals that the learned surface accurately captures both the smile curvature within each maturity slice and the term structure flattening effects as maturity increases. The training converges rapidly, typically requiring 8000 to 10000 total iterations across both optimization stages to reach the specified tolerance thresholds. Calibration quality remains high even in regions with sparse training data, demonstrating effective interpolation guided by the PDE constraints that enforce consistency with the underlying diffusion dynamics. Examination of butterfly spread violations shows that the trained network produces surfaces with positive convexity at 99.9% of evaluated grid points, with the rare violations exhibiting magnitudes below 0.01% of at-the-money variance, far below levels that would enable profitable arbitrage after accounting for transaction costs.

Benchmark comparisons illuminate the advantages of the physics-informed approach relative to alternative methodologies. We implement three competing methods: an unconstrained neural network trained purely on data-fitting objectives without PDE constraints, a parametric SSVI calibration using nonlinear least squares optimization, and cubic spline interpolation with penalty-based smoothing. The unconstrained neural network achieves the lowest in-sample

fitting error with RMSE of 0.2 volatility points, but produces surfaces with substantial arbitrage violations particularly in extrapolation regions beyond the range of training strikes. The butterfly spread condition is violated at approximately 15% of evaluation points with some violations exceeding 1% of ATM variance, creating clear arbitrage opportunities that would be exploited by sophisticated market participants. The parametric SSVI approach guarantees arbitrage-freedom by construction through its constrained parameter space, but exhibits systematic fitting errors with RMSE of 0.8 volatility points due to limited functional flexibility. The model particularly struggles to capture the steep near-term smile, consistently underestimating the curvature for options with maturities below one month. The cubic spline method achieves intermediate performance with RMSE of 0.5 volatility points and predominantly arbitrage-free surfaces, though occasional small violations occur at grid boundaries. However, spline calibration requires extensive manual tuning of knot placement and smoothing parameters, rendering it less attractive for automated real-time applications compared to the PINN approach which achieves superior performance without manual intervention.

## 4.2 Real Market Data Analysis: S&P 500 and VIX Options

Application of the PINN framework to real market data uses end-of-day settlement prices from S&P 500 and VIX index options obtained from publicly available sources covering a six-month period during 2024. The dataset captures diverse market conditions including calm periods with VIX below 15, moderate volatility regimes with VIX between 15 and 25, and elevated volatility episodes with VIX exceeding 25, enabling comprehensive assessment of the methodology's robustness. Each daily snapshot contains 300 to 500 actively traded option contracts across strikes ranging from 50% to 150% of the index level and maturities from one week to one year. Bid-ask spreads are used to construct data-fitting weights in the loss function, with tight spreads below 0.5 volatility points receiving maximum weight and wider spreads above 2.0 points receiving proportionally reduced weight to de-emphasize potentially stale or unreliable quotes. This weighting scheme focuses calibration on liquid benchmark options that represent the most reliable market information.
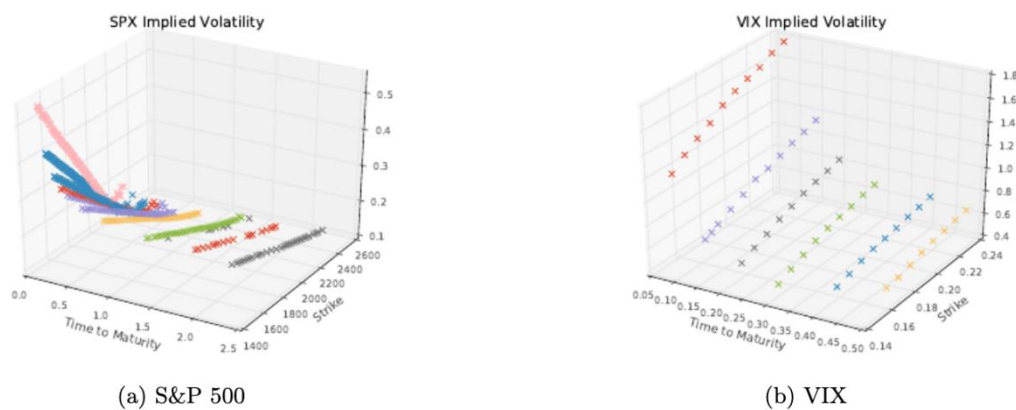


(a) S&P 500                                                                                        (b) VIX

**Figure 3** SPX Implied Volatility and VIX Implied Volatility

The calibration results demonstrate that PINNs successfully fit real market data while maintaining arbitrage-free properties across the vast majority of the surface. As shown in Figure 3, average pricing errors measured in implied volatility terms fall to 0.4-0.6 volatility points for S&P 500 options, well within typical bid-ask spreads and indicating that the network captures market prices as accurately as can be expected given inherent pricing noise and discrete observations. The learned surfaces exhibit realistic features including pronounced put skew for near-term maturities with implied volatilities for 90% moneyness puts exceeding those of at-the-money options by 3-5 volatility points, gradual flattening of the smile for longer-dated options as the skew effect diminishes with increasing maturity, and smooth term structure evolution without artificial oscillations between maturity slices. Arbitrage constraint satisfaction remains high with butterfly spread violations detected at fewer than 1% of evaluation points and calendar spread violations virtually absent. The rare instances of small violations occur in deep out-of-the-money regions where market data is sparse and option delta values fall below 5%, representing strikes that receive minimal trading volume and where pricing is most susceptible to noise.

For VIX options, which present additional modeling challenges due to the mean-reverting nature of volatility indices and their distinct risk characteristics, the PINN framework achieves comparable success with fitting errors of 0.5-0.8 volatility points. The VIX surface displays markedly different smile characteristics compared to S&P 500 options, exhibiting relatively symmetric or even slightly positive skew rather than the negative skew typical of equity options. This reflects the different dynamics of volatility as an asset class, with VIX call options serving as hedges against volatility spikes that often accompany market drawdowns. The joint analysis of S&P 500 and VIX surfaces reveals important insights about market incompleteness and the challenges of consistent multi-asset calibration. While each individual surface satisfies no-arbitrage constraints when considered in isolation, the relationship between the two surfaces exhibits patterns suggesting that market participants employ distinct pricing measures or incorporate different risk preferences when trading volatility derivatives versus equity index options. This observation underscores the practical relevance of incomplete market considerations in real financial applications.

Performance variation across market conditions reveals important patterns regarding the robustness and adaptability of the PINN approach. During calm market periods characterized by VIX below 15 and stable smile shapes, all calibration methods including simple parametric models achieve satisfactory fitting accuracy with errors below 0.5 volatility points. The surfaces exhibit regular smooth patterns well-approximated by standard parameterizations, reducing the advantages of flexible nonparametric approaches. However, during volatile periods when VIX exceeds 20 and smile shapes steepen dramatically, the benefits of physics-informed neural networks become pronounced. Parametric models struggle to capture rapidly changing surface topologies, often requiring parameter bounds or regularization that prevents adequate fitting, resulting in errors exceeding 1.5 volatility points for out-of-the-money options. The PINN framework, by contrast, adapts its functional form to match unusual smile shapes while PDE constraints prevent unrealistic features from emerging, maintaining fitting errors below 0.7 volatility points even during market stress. This robustness to regime changes represents a crucial advantage for practical applications where model reliability during volatile periods is most critical for risk management.

Computational efficiency analysis shows that the PINN implementation achieves calibration times of 45-75 seconds per daily snapshot on standard hardware equipped with modern GPUs, competitive with optimized parametric calibration routines while offering substantially greater flexibility. The bulk of computation time is consumed in the training phase evaluating the composite loss function and its gradients across mini-batches, with each iteration requiring approximately 5-10 milliseconds depending on batch size and network depth. Automatic differentiation enables efficient gradient computation despite the complexity of the loss function incorporating multiple constraint terms, with per-iteration costs scaling linearly in the number of network parameters. Once trained, the neural network provides near-instantaneous volatility surface evaluation at arbitrary strike-maturity points, enabling rapid pricing of exotic derivatives and real-time risk calculations required for dynamic hedging strategies. Forward pass evaluation through the trained network requires only microseconds per query, orders of magnitude faster than traditional methods requiring numerical solution of PDEs or Monte Carlo simulation.

Sensitivity analysis examining the impact of architectural choices and hyperparameter settings provides guidance for practical implementation. Experiments varying network depth from two to six hidden layers reveal that three to four layers offer optimal balance between expressiveness and optimization difficulty, with deeper networks showing diminishing returns in fitting accuracy while increasing training time and occasionally exhibiting instabilities. Hidden layer widths between 50 and 128 neurons prove sufficient for typical volatility surfaces, with larger widths providing marginal benefit for highly complex smiles during volatile market conditions but introducing unnecessary parameters during normal periods. The relative weighting between data-fitting and constraint terms requires problem-specific tuning, with optimal settings depending on data quality, spread magnitudes, and the degree of market incompleteness. We find that final weight allocations placing 40-50% emphasis on data fitting, 25-35% on Dupire constraints, 15-25% on arbitrage conditions, and 5-10% on smoothness regularization work well across most scenarios, though these proportions should be adjusted based on validation performance monitoring.

The practical implications extend beyond academic interest to real-world applications in derivatives trading and risk management. Market makers employing PINN-based surfaces for pricing exotic options benefit from confidence that the underlying volatility structure is arbitrage-free and theoretically consistent, reducing model risk in pricing and hedging activities that could otherwise lead to significant losses when model assumptions are violated. Risk managers utilizing the calibrated surfaces for computing value-at-risk and expected shortfall metrics across large derivative portfolios obtain more reliable risk estimates that properly reflect tail probabilities encoded in the smile shape, improving capital allocation and regulatory compliance. The method's ability to handle incomplete market settings makes it particularly valuable for emerging market equities, commodity derivatives, or other asset classes with limited option liquidity where traditional calibration approaches often fail to produce stable economically sensible results. The transparency afforded by the physics-informed architecture, which explicitly incorporates theoretical constraints rather than functioning as a black box, facilitates model validation and regulatory acceptance in environments requiring explainability.

## 5   CONCLUSION

This paper has developed a comprehensive framework for constructing arbitrage-free volatility surfaces in incomplete markets using Physics-Informed Neural Networks. The proposed methodology successfully integrates empirical market data fitting with theoretical constraints derived from fundamental arbitrage-free pricing principles, producing volatility surfaces that simultaneously achieve high calibration accuracy and satisfy essential no-arbitrage conditions. By incorporating Dupire's local volatility equation, calendar spread constraints, and butterfly spread conditions as differentiable penalty terms within a multi-objective neural network training objective, our approach generates smooth, economically valid implied volatility surfaces that accurately represent complex market patterns including the characteristic smile curvature and term structure evolution observed in equity index options.

The numerical experiments presented demonstrate several key advantages of the PINN-based framework relative to existing parametric and nonparametric methods. First, the approach achieves superior fitting accuracy compared to rigid parametric models while guaranteeing arbitrage-free properties that purely data-driven neural networks fail to satisfy. The calibrated surfaces exhibit root mean squared errors below 0.5 volatility points on both synthetic Heston data and real S&P 500 market observations, well within typical bid-ask spreads, while maintaining butterfly spread violations below 0.1% of at-the-money variance levels across 99% of evaluated points. Second, the methodology exhibits robust

performance across varying market conditions, effectively handling both calm regimes with regular smile patterns and volatile episodes with steep skews, a flexibility that parametric models lack. Third, computational efficiency remains practical for real-world applications, with calibration times of approximately one minute per daily snapshot and near-instantaneous forward evaluation enabling real-time pricing and risk calculations. Fourth, the framework naturally accommodates incomplete market settings where perfect hedging is impossible, producing reasonable volatility surfaces even when theoretical conditions hold only approximately due to market frictions.

The empirical analysis of S&P 500 and VIX option surfaces reveals important insights about market structure and pricing consistency across related derivatives markets. The distinct smile characteristics observed in VIX options compared to equity index options, combined with the challenges of maintaining consistency between the two surfaces, underscore the practical relevance of incomplete market considerations. Market participants appear to employ different pricing principles or incorporate distinct risk preferences when trading volatility derivatives versus equity options, creating subtle inconsistencies that would violate arbitrage-free conditions in a complete market but persist in practice due to hedging constraints and segmentation. The PINN framework's ability to fit each surface independently while respecting individual no-arbitrage constraints provides a pragmatic solution that acknowledges market incompleteness rather than imposing artificial consistency that would degrade fitting quality.

The implications for quantitative finance practice are substantial, as accurate volatility surface models form the foundation of modern derivatives trading and risk management infrastructure. Market makers and proprietary trading desks require surfaces that fit observed prices accurately while remaining free of arbitrage opportunities that would undermine pricing consistency across related instruments. The explicit incorporation of PDE constraints within the PINN architecture provides transparency and interpretability often lacking in black-box machine learning approaches, facilitating model validation procedures required by risk management frameworks and regulatory oversight. The methodology's robustness to data quality variations and ability to produce stable surfaces even with sparse observations addresses a critical practical challenge in markets where option liquidity concentrates in near-the-money strikes and near-term maturities. Extensions to more exotic derivative structures including American options, barrier options, and volatility swaps can leverage the calibrated PINN surfaces, with the arbitrage-free property ensuring that exotic prices remain consistent with liquid vanilla option markets.

Several directions for future research emerge from this work, offering opportunities to extend the framework's capabilities and address remaining limitations. First, incorporating time-varying dynamics to model the evolution of volatility surfaces across consecutive trading days would enable forecasting applications and dynamic recalibration strategies that adapt more rapidly to changing market conditions. The current static framework calibrates each daily surface independently, potentially discarding information contained in the temporal evolution of smile shapes and term structures. Second, investigating alternative PDE constraints beyond Dupire's equation, such as those arising from jump-diffusion models, regime-switching specifications, or rough volatility dynamics, could enhance the method's applicability to markets exhibiting discontinuous price movements or long-memory volatility patterns. The flexibility of the PINN framework permits incorporation of diverse constraint types through appropriate loss function modifications. Third, developing rigorous uncertainty quantification techniques that provide confidence intervals or posterior distributions over the predicted volatility surface would support risk-aware decision making and model risk assessment. Current point estimates do not capture the uncertainty arising from limited data and model approximation errors. Fourth, extending the methodology to multi-asset settings where correlations between underlying assets introduce additional complexity offers promising avenues for applications in portfolio risk management and cross-asset derivatives pricing.

The integration of physics-informed approaches with modern machine learning architectures represents a powerful paradigm for addressing complex financial modeling challenges. By combining the flexibility and computational efficiency of neural networks with the structure and domain knowledge encoded in financial theory through governing PDEs, researchers and practitioners can develop models that are both empirically accurate and theoretically grounded. The success of PINNs in volatility surface construction suggests that similar frameworks may prove valuable for other financial problems involving differential equations or optimization under constraints, including interest rate curve modeling, credit risk assessment, and optimal execution strategies. As machine learning continues to transform quantitative finance, approaches that respect domain knowledge while harnessing computational power will play an increasingly central role in advancing the field's capabilities and maintaining the reliability required for high-stakes financial applications.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Bae HO, Kang S, Lee M. Option Pricing and Local Volatility Surface by Physics-Informed Neural Network. Computational Economics, 2024, 64(5): 3143-3159.

[2] Ma Z, Chen X, Sun T, et al. Blockchain-based zero-trust supply chain security integrated with deep reinforcement learning for inventory optimization. Future Internet, 2024, 16(5): 163.

[3] Almeida C, Freire G. Pricing of index options in incomplete markets. Journal of Financial Economics, 2022, 144(1): 174-205.

[4] Sun T, Yang J, Li J, et al. Enhancing auto insurance risk evaluation with transformer and SHAP. IEEE Access, 2024.

[5] Ruf J, Wang W. Neural networks for option pricing and hedging: a literature review. arXiv preprint arXiv:1911.05620, 2019.

[6] Cao Y, Liu X, Zhai J. Option valuation under no-arbitrage constraints with neural networks. European Journal of Operational Research, 2021, 293(1): 361-374.

[7] Berner J, Grohs P, Jentzen A. Neural Networks can detect model-free static arbitrage strategies. arXiv preprint, 2024.

[8] Raissi M, Perdikaris P, Karniadakis GE. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. Journal of Computational Physics, 2019, 378: 686-707.

[9] Dhiman A, Hu Y. Physics Informed Neural Network for Option Pricing. arXiv preprint arXiv:2312.06711, 2023.

[10] Ge Y, Wang Y, Liu J, et al. GAN-Enhanced Implied Volatility Surface Reconstruction for Option Pricing Error Mitigation. IEEE Access, 2025.

[11] Zhang W, Li L, Zhang G. A two-step framework for arbitrage-free prediction of the implied volatility surface. Quantitative Finance, 2023, 23(1): 21-34.

[12] Zheng Y, Yang Y, Chen B. Incorporating prior financial domain knowledge into neural networks for implied volatility surface prediction. Proceedings of the 27th ACM SIGKDD Conference, 2021: 3968-3975.

[13] Nyah E E, Onwuka D O, Arimanwa J I, et al. Adaptive neuro-fuzzy inference system optimization of natural rubber latex modified concrete's mechanical Properties. Scientific Reports, 2025, 15(1): 20624.

[14] Wang S. Arbitrage-free neural-SDE market models of traded options. Doctoral dissertation, University of Oxford, 2022.

[15] Mai N T, Cao W, Liu W. Interpretable knowledge tracing via transformer-Bayesian hybrid networks: Learning temporal dependencies and causal structures in educational data. Applied Sciences, 2025, 15(17): 9605.

[16] Cao W, Mai N T, Liu W. Adaptive knowledge assessment via symmetric hierarchical Bayesian neural networks with graph symmetry-aware concept dependencies. Symmetry, 2025, 17(8): 1332.

[17] Chen S, Liu Y, Zhang Q, et al. Multi-Distance Spatial-Temporal Graph Neural Network for Anomaly Detection in Blockchain Transactions. Advanced Intelligent Systems, 2025: 2400898.

[18] Mai N T, Cao W, Wang Y. The global belonging support framework: Enhancing equity and access for international graduate students. Journal of International Students, 2025, 15(9): 141-160.

[19] Chen S, Liu Y, Zhang Q, et al. Multi-Distance Spatial-Temporal Graph Neural Network for Anomaly Detection in Blockchain Transactions. Advanced Intelligent Systems, 2025: 2400898.

[20] Zhang Q, Chen S, Liu W. Balanced Knowledge Transfer in MTTL-ClinicalBERT: A Symmetrical Multi-Task Learning Framework for Clinical Text Classification. Symmetry, 2025, 17(6): 823.

[21] Ren S, Jin J, Niu G, et al. ARCS: Adaptive Reinforcement Learning Framework for Automated Cybersecurity Incident Response Strategy Optimization. Applied Sciences, 2025, 15(2): 951.

[22] Tan Y, Wu B, Cao J, et al. LLaMA-UTP: Knowledge-Guided Expert Mixture for Analyzing Uncertain Tax Positions. IEEE Access, 2025.

[23] Zheng W, Liu W. Symmetry-Aware Transformers for Asymmetric Causal Discovery in Financial Time Series. Symmetry, 2025, 17(10): 1591.

[24] Liu Y, Ren S, Wang X, et al. Temporal logical attention network for log-based anomaly detection in distributed systems. Sensors, 2024, 24(24): 7949.

[25] Hu X, Zhao X, Wang J, et al. Information-theoretic multi-scale geometric pre-training for enhanced molecular property prediction. PLoS One, 2025, 20(10): e0332640.

[26] Zhang H, Ge Y, Zhao X, et al. Hierarchical deep reinforcement learning for multi-objective integrated circuit physical layout optimization with congestion-aware reward shaping. IEEE Access, 2025.

[27] Wang J, Zhang H, Wu B, et al. Symmetry-Guided Electric Vehicles Energy Consumption Optimization Based on Driver Behavior and Environmental Factors: A Reinforcement Learning Approach. Symmetry, 2025, 17(6): 930.

[28] Han X, Yang Y, Chen J, et al. Symmetry-Aware Credit Risk Modeling: A Deep Learning Framework Exploiting Financial Data Balance and Invariance. Symmetry, 2025, 17(3): 20738994.

[29] Hu X, Zhao X, Liu W. Hierarchical Sensing Framework for Polymer Degradation Monitoring: A Physics-Constrained Reinforcement Learning Framework for Programmable Material Discovery. Sensors, 2025, 25(14): 4479.

[30] Qiu L. Reinforcement Learning Approaches for Intelligent Control of Smart Building Energy Systems with Real-Time Adaptation to Occupant Behavior and Weather Conditions. Journal of Computing and Electronic Information Management, 2025, 18(2): 32-37.

# ANALYSIS OF THE IMPACT OF SHOOTING SPORTS ON THE PSYCHOLOGICAL QUALITY OF TEENAGERS——BASED ON THE INVESTIGATION AND ANALYSIS OF SHOOTING CLUBS

YunYi Wang

*Zhixin High School, Guangzhou 510060, Guangdong, China.*
*Corresponding Email: usk19231899@outlook.com*

**Abstract:** The requirements for teenagers' psychological resilience are getting higher in modern society, shooting sports are beneficial to teenagers' mental health. This paper studies the impact of shooting on teenagers' psychological factors such as concentration, emotional control ability, adjustment ability, endurance and stress adaptation ability. By means of the cluster sampling, a questionnaire survey is conducted in shooting clubs, and the results of the questionnaire survey are systematically analyzed. The analysis shows that through shooting training, the three comprehensive indicators of self-control, concentration and psychological stability of the trainees have been significantly improved, and the effect is more obvious for trainees who have trained for more than one year. Moreover, the surveyed trainees all have a positive willingness to promote shooting sports.
**Keywords:** Shooting sports; Psychological quality; Sampling survey

## 1 INTRODUCTION

Shooting sports are gradually becoming a popular sports activity among teenagers. It originates from hunting and military activities and is a new sport added to human activities in the era of hot weapons. It is also the most direct inheritance and embodiment of modern military skills. In modern society, the requirements for teenagers' psychological resilience are getting higher, and a high psychological quality is the guarantee for learning and work. Shooting sports can cultivate teenagers' concentration, emotional control ability, and psychological resilience, and play an important role in enhancing teenagers' self-confidence and self-identity.

In order to study the impact of shooting sports on the psychological education of teenagers, I conducted a questionnaire survey from July to August 2025 during the summer vacation. I designed a questionnaire and conducte the survey among members of Guangzhou Shuimu Shooting Club and Haizhifeng Youth Shooting Training Base. The respondents are mainly people aged 25 and below.

## 2 OBJECTIVES AND METHODS

### 2.1 Survey Participants

This survey is conducted by using cluster sampling, selecting two shooting clubs in Guangzhou and Hefei of Anhui Province for the investigation. The survey subjects are the trainees who participate in shooting training, and there are no restrictions on the gender and age of the trainees. In cases involving younger participants, the questionnaires are completed with the assistance of their guardians.

### 2.2 Survey Tools

This survey is conducted in the form of a questionnaire, which is divided into three parts: basic information, psychological condition rating scale and recommendation intention scale. The basic information covered basic questions such as age, gender, current educational stage, and duration of shooting training.

The psychological condition rating scale is consisted of 19 questions, covering five aspects: concentration, self-control, stress resistance, psychological stability, and comprehensive psychological quality. Each question is presented in a Likert scale, with five options ranging sequentially from "completely not conforms", "less conforms", "moderately conforms", "basically conforms", to "completely conforms". All items are positive questions. After verification. the Cronbach's alpha coefficient for the scale questions is 0.959, indicating a high level of reliability.

The recommendation intention scale consists of a single item, and is used to reflect the trainees' willingness to recommend shooting sports. It also employs a Likert-scale with five predefined response options.

### 2.3 Statistical Analysis Methods

This questionnaire analysis is conducted using Python 3.11.5 and SPSS 27.0, with the method of mean comparison and variance analysis; the analysis of influencing factors adopts binary categorical variable Logistic regression.

## 3 ANALYSIS OF BASIC SITUATION

### 3.1 Distribution of School Age of the Survey Participants

A total of 109 questionnaires are actually collected in this survey, of which 2 respondents have never received shooting training, so the actual number of questionnaires include in the analysis is 107. According to the age range of the respondents, they are divided into 6 age groups.

**Table 1** Age Distribution

| Age | Frequency | Proportion |
|---|---|---|
| Preschool | 4 | 3.74% |
| Primary School | 68 | 63.55% |
| Junior High School | 17 | 15.89% |
| High School | 6 | 5.61% |
| Others | 12 | 11.21% |
| Total | 107 | 100% |

By age group, students in the primary school age range of 6-11 years, account for the highest proportion of 57.01%; followed by the junior high school age group of 12-14 years with a proportion of 25.23%. The combined proportion of these two age groups reaches 82.24%, indicating that the primary and junior high school age groups are the main groups practicing shooting.

### 3.2 Gender Distribution of the Survey Participants

From the perspective of gender, 67 of the respondents are male and 40 are female, accounting for 62.62% and 37.38% respectively, indicating that the shooters are mainly male.

**Table 2** Distribution of Gender

| Gender | Frequency | Proportion |
|---|---|---|
| Male | 67 | 62.62% |
| Female | 40 | 37.38% |
| Total | 107 | 100% |

### 3.3 Distribution of Training Duration of the Survey Participants

All 107 individuals surveyed have practiced shooting, however, their training duration varied significantly. Among them, 65 individuals, accounting for 60.75% of the total, have practiced for less than 6 months, which is the highest proportion. Additionally, 18.69% and 12.15% of the participants have practiced for less than one year and less than three years respectively. The participants with training durations of more than five years are relatively few, accounting for only 8.41%, and they are aged 18 or above, with the oldest being 26 years old. This indicates that the shooters are mainly new learners, and relatively few have persisted in learning for a long period.

**Table 3** Distribution of Practice Duration

| Options | Frequency | proportion |
|---|---|---|
| Less than 6 months | 65 | 60.75% |
| 6 months to 1 year | 20 | 18.69% |
| 1 to 3 years | 13 | 12.15% |

| | | | |
|---|---|---|---|
| 3 to 5 years | 0 | | 0% |
| More than 5 years | 9 | | 8.41% |
| Has not participated in any shooting training | 0 | | 0% |
| Total | 107 | | 100% |

## 3.4 Validation of the Survey Scale's Reliability and Validity

The reliability of 19 scale questions is tested, yielding a Cronbach's alpha coefficient of 0.959, which exceeds 0.7. This result indicates that the questionnaire possesses a high level of reliability.

The validity of the questionnaire is tested using the principal component analysis method, with variance maximization rotation. According to the purpose of questionnaire design, 5 principal components are extracted. After testing, there are two questions, namely "competition pressure management" and "high-pressure decision-making ability" among the 19 scale questions, whose matrix loadings on each factor are all less than 0.5, indicating that they are invalid questions. In the subsequent analysis, these two questions are excluded, and the remaining 17 questions are discussed and analyzed.

## 4 RESULT ANALYSIS

### 4.1 Factor Analysis

The data is first subjected to KMO and Bartlett tests for the 17 questions of the questionnaire scale. The results show that the sample data is suitable for factor analysis, with a KMO value of 0.935 and a Bartlett's test p-value of less than 0.01.

Combining the five aspects of the questionnaire design, the principal component analysis is used to extract the principal factors, and the factor loadings are rotated using the variance maximization method. The 5 principal factors, based on their factor loadings, account for 75.79% of the total variance, and the rotated component matrix is presented in the table below.

**Table 4** Rotated Component Matrix

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Daily concentration ability | .761 | | | | |
| Long-term psychological resilience | .692 | | | | |
| Emotional transfer | .610 | | | | |
| Reduction in impulsive behavior | .576 | | | | |
| Impact of external evaluation | .531 | | | | |
| Muscle Relaxation Techniques | .522 | | | | |
| Heart rate control | | .771 | | | |
| Recovery of emotions after a mistake | | .684 | | | |
| Operational stability | | .628 | | | |
| Task switching and restoration | | .627 | | | |
| Visual focus | | | .784 | | |
| Resist environmental interference | | | .688 | | |
| Execution of the training plan | | | | .758 | |
| Suppressing impulsive shooting | | | | .667 | |
| respiratory control | | | | .515 | |
| Long-term concentration endurance | | | | | .779 |
| Stress adaptation | | | | | .664 |

Based on the factor loadings, factor 1 can be named as the comprehensive psychological quality factor, factor 2 as the psychological stability factor, factor 3 as the concentration ability factor, factor 4 as the self-control ability factor, and factor 5 as the endurance and stress adaptation factor.

### 4.2 Gender Differences in Participants

The results of the variance analysis of the five factors based on gender as the classification variable are shown in Table 5. The results indicate that, at the 0.05 significance level, here are no significant gender differences in four out of the five factors, however, there are significant differences in the factor of anti-interference ability. Further analysis reveals that the mean value of the anti-interference ability factor for females is significantly higher than that for males, indicating that shooting training effectively enhances the anti-interference ability of females.

**Table 5** Analysis of Variance of the Impact of Gender Factors on the Factors

| | | | sum of squares | DOF | mean square | F | significance |
|---|---|---|---|---|---|---|---|
| Factor 1 (Comprehensive Psychological Quality) * Gender | Group Comparison | (combination) | .060 | 1 | .060 | .060 | .807 |
| | Intra-group | | 105.940 | 105 | 1.009 | | |
| | Total | | 106.000 | 106 | | | |
| Factor 2 (Psychological Stability) * Gender | Group Comparison | （combination） | .383 | 1 | .383 | .381 | .538 |
| | Intra-group | | 105.617 | 105 | 1.006 | | |
| | Total | | 106.000 | 106 | | | |
| Factor 3 (Concentration) * Gender | Group Comparison | （combination） | 6.206 | 1 | 6.206 | 6.530 | .012 |
| | Intra-group | | 99.794 | 105 | .950 | | |
| | Total | | 106.000 | 106 | | | |
| Factor 4 (Self-control) * Gender | Group Comparison | （combination） | .538 | 1 | .538 | .535 | .466 |
| | Intra-group | | 105.462 | 105 | 1.004 | | |
| | Total | | 106.000 | 106 | | | |
| Factor 5 (Endurance and Stress Adaptation) * Gender | Group Comparison | （combination） | 1.587 | 1 | 1.587 | 1.596 | .209 |
| | Intra-group | | 104.413 | 105 | .994 | | |
| | Total | | 106.000 | 106 | | | |

## 4.3 The Impact of the Educational Stage on Trainees' Psychology

The participants surveyed cover five stages: preschool, primary school, junior high school, senior high school, and other. The results of the variance analysis indicate that there are significant differences in the impact of shooting training on the execution factor across different educational stages at the 0.1 significance level, while the impact on other factors is not significant.

**Table 6** Analysis of Variance of the Impact of Educational Stage on Factors

| | | Sum of Squares | DOF | Mean Square | F | Significance |
|---|---|---|---|---|---|---|
| Factor 1 (Comprehensive Psychological Quality) * Educational stage | Group Comparison | 4.519 | 4 | 1.130 | 1.136 | .344 |
| | Intra-group | 101.481 | 102 | .995 | | |
| | Total | 106.000 | 106 | | | |
| Factor 2 (Psychological Stability) * Educational Stage | Group Comparison | 2.390 | 4 | .597 | .588 | .672 |
| | Intra-group | 103.610 | 102 | 1.016 | | |
| | Total | 106.000 | 106 | | | |
| Factor 3 (Concentration) * Educational Stage | Group Comparison | 5.945 | 4 | 1.486 | 1.515 | .203 |
| | Intra-group | 100.055 | 102 | .981 | | |
| | Total | 106.000 | 106 | | | |
| Factor 4 (Self-control) * Educational stage | Group Comparison | 7.847 | 4 | 1.962 | 2.039 | .095 |
| | Intra-group | 98.153 | 102 | .962 | | |
| | Total | 106.000 | 106 | | | |
| Factor 5 (Endurance and Stress Adaptation) * Educational Stage | Group Comparison | 3.408 | 4 | .852 | .847 | .499 |
| | Intra-group | 102.592 | 102 | 1.006 | | |
| | Total | 106.000 | 106 | | | |

Considering that the educational stage is a multiclass variable, multiple comparisons are adopted to further investigate the impact of different educational stages on the self-control ability of the trainees. The statistical quantity calculated by the double comparison is LSD, and the formula is the following:

$$LSD = t_{\alpha/2} \sqrt{MSE(\frac{1}{n_i} + \frac{1}{n_j})} \tag{1}$$

Based on the actual situation of this survey, the value of t in the formula is taken as the critical value of a two-sided test with a degree of freedom of 102, and the significance level is set at 0.05; MSE is the within-group mean square error, which is 0.962285; "ni" is the sample size of the i-th educational stage. The absolute value of the mean difference between the i-th and j-th educational stages is calculated, and compared with the corresponding LSD value. If the former is less than the latter, it indicates that there is no significant difference between the two educational stages; if the former is greater than the latter, it indicates that there is a significant difference between the two educational stages.

The test results show that significant differences exist between educational stage 1 and stage 5, and between stage 2 and stage 5, while the other educational stages are not obvious. Since the mean score of the factor in educational stage 5 is significantly higher than that of stages 1 and 2, it indicates that the students in the "other" educational stage have a significantly greater impact on shooting in terms of self-control compared to those in the preschool and primary school stages.

Further analysis of the sample data reveals that there are a total of 12 students in educational stage 5, all of whom are

adults (aged 18 or above), and 5 of them are over 25 years old.

## 4.4 The Impact of Training Duration on the Mind

From the perspective of training duration, the results of the variance analysis indicate that there are significant differences in the effects of practice duration on aspects such as psychological stability, anti-interference ability, and execution ability at the 0.1 significance level. Specifically, the psychological stability is the best during the 1-3 year practice period; the anti-interference ability of the trainees with a practice duration of 5 years or more is the strongest, followed by those with a training duration of 1-3 years, who are significantly higher than those with other training durations; the execution ability of the trainees with a practice duration of 5 years or more is the best.

**Table 7** Analysis of Variance of the Psychological Impact Caused by Practice Duration

|  |  |  | Sum of Squares | DOF | Mean Square | F | Significance |
|---|---|---|---|---|---|---|---|
| Factor 1 (Comprehensive Psychological Quality) * Practice Duration | Group Comparison | (combination) | 2.230 | 3 | .743 | .738 | .532 |
|  | Intra-group |  |  | 103 | 1.007 |  |  |
|  | Total |  |  | 106 |  |  |  |
| Factor 2 (Psychological Stability) * Practice Duration | Group Comparison | (combination) | 11.468 | 3 | 3.823 | 4.165 | .008 |
|  | Intra-group |  |  | 103 | .918 |  |  |
|  | Total |  |  | 106 |  |  |  |
| Factor 3 (Concentration) * Practice Duration | Group Comparison | (combination) | 12.714 | 3 | 4.238 | 4.679 | .004 |
|  | Intra-group |  |  | 103 | .906 |  |  |
|  | Total |  |  | 106 |  |  |  |
| Factor 4 (Self-control) * Practice duration | Group Comparison | (combination) | 6.621 | 3 | 2.207 | 2.287 | .083 |
|  | Intra-group |  |  | 103 | .965 |  |  |
|  | Total |  |  | 106 |  |  |  |
| Factor 5 (Endurance and Stress Adaptation) * Duration of Training | Group Comparison | (combination) | 2.772 | 3 | .924 | .922 | .433 |
|  | Intra-group |  |  | 103 | 1.002 |  |  |
|  | Total |  |  | 106 |  |  |  |

## 5 REGRESSION ANALYSIS

### 5.1 Multivariate Linear Regression Analysis

In order to further verify the causal relationship between shooting sports and psychological qualities, this study constructs the following regression model:

$$fa_i = \alpha_0 + \alpha_1 ge_i + \alpha_2 ed_i + ex_i + \mu_i \tag{2}$$

In the above equation, the value of i ranges from 1 to 5. fa represents the aforementioned factors, $\alpha_0$ is a constant term, ge represents gender, ed represents the educational stage (indicating age), ex represents the training duration, and $\mu_i$ is a random term. Using the OLS regression method, the regression results of the model with these 5 factors as dependent variables are as follows:

**Table 8** Regression Results of Comprehensive Psychological Quality

|  | Coef | Std err | t | P > \|t\| |
|---|---|---|---|---|
| Gender | -0.0467 | 0.207 | -0.226 | 0.822 |
| Educational Stage | 0.1499 | 0.122 | 1.233 | 0.220 |
| Training Duration | 0.0292 | 0.104 | 0.280 | 0.780 |
| Const | -0.3724 | 0.337 | -1.106 | 0.271 |

As shown in Table 8, the factors of gender, educational stage, and duration of shooting practice have no significant impact on Factor 1 (Comprehensive Psychological Quality). The Comprehensive psychological quality includes six indicators such as daily concentration, long-term psychological resilience, emotion transfer, reduction of impulsive behavior, influence of external evaluation, and muscle relaxation techniques. No significant correlation is observed between the duration of shooting training and this composite indicator.

**Table 9** Regression Results of Psychological Stability

|  | Coef | Std err | t | P > \|t\| |
|---|---|---|---|---|
| Gender | -0.1377 | 0.202 | -0.683 | 0.496 |
| Educational Stage | -0.1644 | 0.119 | -1.403 | 0.160 |
| Training Duration | 0.2903 | 0.102 | 2.854 | 0.005 |

| | | |
|---|---|---|
| Const | 0.1040 0.329 0.317 0.752 | |

Factor 2 is characterized by psychological stability, which is composed of four indicators: heart rate control, emotional recovery after mistakes, operational stability, and task-switching recovery. From Table 9, it can be seen that gender and educational stage have no significant impact on psychological stability, while the duration of shooting training has a significant correlation with psychological stability. This indicates that as the duration of shooting practice increases, the psychological stability of the trainees improves. Thus, it can be concluded that the duration of shooting training has a positive impact on psychological stability.

**Table 10** Regression Results of Concentration Ability

| | Coef | Std err | t | P > \|t\| |
|---|---|---|---|---|
| Gender | 0.4648 | 0.193 | 2.404 | 0.018 |
| Educational Stage | -0.1510 | 0.114 | -1.328 | 0.187 |
| Training Duration | 0.3179 | 0.098 | 3.259 | 0.002 |
| Const | -0.8119 | 0.315 | -2.576 | 0.011 |

Factor 3, Concentration, is consisted of two indicators: visual focus and resistance to environmental interference. As can be seen from the results in Table 10, there is a significant positive correlation between gender, practice duration and concentration ability. This indicates that after shooting training, the concentration ability of girls has improved more significantly than that of boys; and the longer the shooting practice time for teenagers, the stronger their concentration ability. It can be seen that shooting sports have a significant enhancing effect on concentration ability.

**Table 11**    Regression Results of Self-Control Ability

| | Coef | Std err | t | P > \|t\| |
|---|---|---|---|---|
| Gender | 0.0206 | 0.203 | 0.101 | 0.919 |
| Educational Stage | 0.1619 | 0.119 | 1.356 | 0.178 |
| Training Duration | -0.0916 | 0.102 | 0.895 | 0.373 |
| Const | -0.6061 | 0.331 | -1.833 | 0.070 |

Factor 4, Self-control, is composed of four indicators: execution of the training plan, suppression of impulsive shooting, and adjustment of breathing. From the results in Table 11, the influence of gender, educational stage, and practice duration on self-control is not significant. This might be due to the limited sample size, however, the results may differ if the sample size is expanded.

**Table 12**    Regression Results of Endurance and Stress Adaptation

| | Coef | Std errt | | P > \|t\| |
|---|---|---|---|---|
| Gender | 0.1682 | 0.205 | 0.820 | 0.414 |
| Educational Stage | 0.2197 | 0.121 | 1.821 | 0.071 |
| Training Duration | -0.1098 | 0.104 | -1.061 | 0.291 |
| Const | -0.6018 | 0.334 | -1.800 | 0.075 |

Factor 5, Endurance and Stress Adaptation, is consisted of two indicators: long-term concentration endurance and stress adaptation. From the regression results in Table 12, it can be seen that the educational stage has the most significant impact on this comprehensive indicator. In the questionnaire, the educational stage is divided into five stages: preschool, primary school, junior high school, senior high school, and others. It can be seen that the educational stage is an indicator representing age. The educational stage has a positive impact on endurance and stress adaptation. For the surveyed adolescent trainees, their endurance and stress adaptation abilities become stronger as age increases.

## 5.2 Logistic Regression Analysis of Recommendation Intention

The five levels of recommendation intention are converted into binary values in this article. The two options of "basically conform" and "completely conform" are defined as 1; the three options of "completely not conform", "less conform", and "moderately consistent" are defined as 0. A binary Logistic regression analysis is conducted to examine the trainees' recommendation intentions for shooting sports, with the five factors serving as independent variables.

**Table 13** Logistic Regression Results of Recommendation Intent

| | Coef | Std err | z | P |
|---|---|---|---|---|
| Factor 1 (Comprehensive Psychological Quality) | 2.437 | 0.544 | 4.480 | 0.000 |
| Factor 2 (Psychological Stability) | 0.547 | 0.313 | 1.747 | 0.081 |
| Factor 3 (Concentration) | 0.565 | 0.342 | 1.652 | 0.098 |
| Factor 4 (Self-control) | 0.961 | 0.331 | 2.903 | 0.004 |
| Factor 5 (Endurance and Stress Adaptation) | 0.766 | 0.338 | 2.265 | 0.024 |
| Const | 1.966 | .436 | 4.512 | 0.000 |

The regression results in Table 13 show that at all five factors significantly influence the willingness of trainees to recommend at the 0.1 significance level. Among these factors, the odd ratio (OR) of the comprehensive psychological quality factor is the highest at 11.435, indicating its strongest impact on the willingness to recommend; the second highest is the self-control, endurance and stress adaptation factor, with ORs of 2.615 and 2.15 respectively; the ORs of the psychological stability and concentration factors are 1.727 and 1.759 respectively. In summary, the willingness of trainees to recommend is influenced by all five factors, but the most significant is the comprehensive psychological quality factor, indicating that trainees value the impact of shooting on the improvement of their comprehensive psychological quality the most, which is the main factor for trainees to consider when recommending.

## 6 CONCLUSION

The vast majority of adolescents are in full-time education, and primary and secondary school students are under the pressure of further education, leading to frequent mental health issues. Mental health issues have become an important factor affecting students' learning and life. This article aims to study the positive impact of shooting sports on adolescent mental health. As can be seen from the above analysis, shooting sports have a significant impact on adolescent mental health, and different ages, genders, and training durations play different roles in this process.

In terms of the impact of shooting training on mental health, the age factor is mainly reflected in the gradual increase in stress adaptation ability, or stress resistance, and endurance of the surveyed trainees as they grow older. In terms of gender, it is mainly reflected in the significant improvement in concentration of female trainees through shooting training. This indicates that compared to males, females are more easily influenced by their surroundings, and shooting training can improve their ability to resist interference and focus.

The duration of shooting training is a crucial factor in enhancing the mental health of adolescents. Research indicates that the longer the shooting training duration, the higher the individual's psychological resilience and emotional stability. Analysis reveals that through shooting training, students significantly improve in three key areas: self-focus, psychological stability, and their ability to endure and adapt to stress. For students who have trained for over a year, the effects are even more pronounced.

In terms of willingness to recommend shooting sports, the five indicators of comprehensive psychological quality, psychological stability, concentration, self-control, endurance, and stress adaptation are all significant. From the perspective of improving comprehensive psychological quality, the willingness to recommend shooting sports is the highest, followed by the willingness to recommend shooting sports from the perspective of improving self-control, and then endurance and stress adaptation, psychological stability, and concentration.The findings suggest that improvement in any single dimension of psychological quality is sufficient to foster a strong willingness in trainees to promote shooting sports.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Huang Gang. Experimental study on the influence of mindfulness training on mental fatigue of young shooters. Liaoning Sport Sciencs and Technology, 2023(3): 114-119.
[2] Cashel M L, Brooks S, Crawford T P, et al. Diversity training in clinical child and pediatric psychology programs: Results from a survey of training directors. Training and Education in Professional Psychology, 2025, 19(3): 217–225.
[3] Fujiwara K, Varghese V, Chikaraishi M, et al. Does response lag affect travelers' stated preference? Evidence from a real-time stated adaptation survey. Transportation, 2025, 52: 693–713.
[4] Puspitasari S, Varghese V, Chikaraishi M, et al. Exploring the effects of congestion pricing on travel behavior responses using real-time, context-aware, stated-preference data. Journal of the Eastern Asia Society for Transportation Studies, 2021, 14: 199–214.
[5] Ghaith Ahmad, Ma Huimin, Ding Xiuhao. Investigating user acceptance of fully autonomous vehicles in MENA region. Transportation Research Part F: Traffic Psychology and Behaviour, 2025, 114: 1024-1041.
[6] Keigo Akimoto, Fuminori Sano, Junichiro Oda. Impacts of ride and car-sharing associated with fully autonomous cars on global energy consumptions and carbon dioxide emissions. Technological Forecasting and Social Change, 2022, 174: 121311.
[7] Ormanci U, Cepni S. The Effect of Web-Assisted Guided Inquiry Approach on Students' Systems Thinking Skills. Journal of Science Education and Technology, 2025: 1-25.
[8] Xiaoping Tang, Zijun Shen, Muhammad Idrees Khan, et al. A sociological investigation of the effect of cell phone use on students' academic, psychological, and socio-psychological performance. Frontiers in Psychology, 2025, 16: 161474340-1474340.
[9] Rebecca F Slykerman, Eileen Li, Edwin A Mitchell. Students' Experience of Online University Education During the COVID-19 Pandemic: Relationships to Psychological Health. Student Success, 2022, 13(1): 32-40.

[10] Frolova E V, Rogach O V, Razov P V. The effect of the non-financial stimulation system on psychological health of higher school workers. Problemy sotsial'noi gigieny, zdravookhraneniia i istorii meditsiny, 2023, 31(2): 284-289.

[11] Liangqun Y, Jian L. Quantitative Evaluation Method of Psychological Quality of College Teachers Based on Fuzzy Logic. International Journal of Information Technology and Web Engineering (IJITWE), 2024, 19(1): 1-19.

[12] Bao Y, Wang J, Huang H, et al. Preliminary investigation of the association between air pollution exposure and childhood asthma hospitalizations from 2015 to 2018 in East China. Frontiers in Public Health, 2025, 13: 1527214-1527214.

# HIERARCHICAL GNN FRAMEWORK FOR ENERGY-AWARE SCHEDULING IN GPU-ACCELERATED DISTRIBUTED SYSTEMS

Isabella Marino, Lukas Hoffmann*
*Department of Computer Science, Technical University of Munich (TUM), Germany.*
*Corresponding Author: Lukas Hoffmann, Email: lukas.hoffmann@gmail.com*

**Abstract:** Graphics Processing Units have become indispensable computational accelerators in modern distributed computing systems, powering applications ranging from deep learning to scientific simulations. However, the increasing computational demands and energy consumption of GPU-accelerated systems pose significant challenges for resource management and scheduling. Traditional scheduling algorithms often fail to capture the complex hierarchical structure and dynamic dependencies inherent in distributed GPU environments, leading to suboptimal energy efficiency and performance degradation. This paper proposes a novel hierarchical Graph Neural Network framework for energy-aware scheduling in GPU-accelerated distributed systems. The framework leverages the representational power of GNNs to model the complex interactions between computational tasks, GPU resources, and energy constraints at multiple hierarchical levels. By incorporating graph-structured representations of workload dependencies, resource topologies, and energy profiles, the proposed framework enables adaptive scheduling decisions that jointly optimize task completion time and energy consumption. Experimental results demonstrate that the proposed approach achieves up to 36 percent reduction in energy consumption while maintaining quality of service requirements compared to state-of-the-art scheduling methods. The hierarchical architecture effectively captures both fine-grained GPU-level characteristics and coarse-grained cluster-level dynamics, enabling scalable and efficient scheduling for large-scale distributed systems.
**Keywords:** Graph neural networks; Energy-aware scheduling; GPU computing; Distributed systems; Hierarchical framework; Resource management; Deep learning

## 1 INTRODUCTION

The rapid advancement of artificial intelligence and data-intensive computing has driven unprecedented demand for high-performance computational resources, with Graphics Processing Units emerging as the cornerstone of modern computing infrastructure. GPU-accelerated systems have demonstrated remarkable capabilities in accelerating complex workloads including deep neural network training, scientific simulations, and large-scale data analytics. Major technology companies and research institutions have deployed extensive GPU clusters comprising thousands of interconnected devices to support their computational requirements [1]. However, this explosive growth in GPU deployment has introduced critical challenges in energy management and resource scheduling that demand innovative solutions.

Energy consumption has become a primary concern in large-scale GPU-accelerated distributed systems due to both economic and environmental considerations. Modern data centers consume enormous amounts of electrical power, with GPU clusters contributing substantially to the overall energy footprint [2]. The operational costs associated with energy consumption and cooling infrastructure can account for a significant portion of the total cost of ownership for computing facilities. Furthermore, the environmental impact of energy-intensive computing has prompted increased scrutiny from regulatory bodies and heightened awareness within the research community regarding sustainable computing practices [3]. These factors collectively underscore the urgent need for energy-efficient scheduling strategies that can reduce power consumption without compromising computational performance or quality of service.

Traditional scheduling algorithms designed for CPU-based systems often prove inadequate when applied to GPU-accelerated distributed environments due to fundamental architectural differences and unique operational characteristics. GPUs exhibit distinct power consumption patterns influenced by multiple factors including workload characteristics, memory access patterns, data transfer overheads, and dynamic voltage and frequency scaling capabilities [4]. The performance saturation characteristics of GPUs differ significantly from CPUs, particularly for computational workloads with varying matrix sizes and parallelism degrees. Understanding these performance differences is crucial for effective scheduling decisions that maximize throughput while minimizing energy consumption. The hierarchical nature of distributed GPU systems, encompassing individual GPU cores, multi-GPU nodes, and cluster-level interconnections, introduces additional complexity that conventional flat scheduling approaches struggle to address effectively [5].

Graph Neural Networks have emerged as a powerful paradigm for learning representations from graph-structured data, demonstrating exceptional performance across diverse application domains including social network analysis, molecular property prediction, and combinatorial optimization [6]. The fundamental strength of GNNs lies in their ability to capture complex relational patterns and dependencies through iterative message passing and aggregation

mechanisms. Recent advances in GNN architectures, including attention mechanisms, hierarchical pooling, and adaptive sampling strategies, have significantly enhanced their scalability and expressiveness [7]. These developments have motivated researchers to explore GNN applications in system optimization and resource management problems where underlying structures can be naturally represented as graphs.

The scheduling problem in GPU-accelerated distributed systems exhibits inherent graph structure at multiple levels of abstraction. At the task level, computational workloads can be represented as directed acyclic graphs capturing data dependencies and execution precedence constraints. At the resource level, GPU clusters form complex topologies with heterogeneous devices interconnected through various communication fabrics. At the system level, energy profiles, thermal characteristics, and performance metrics constitute additional graph-structured relationships that influence scheduling decisions [8]. This multi-level graph structure provides strong motivation for developing GNN-based scheduling frameworks that can effectively leverage these inherent patterns to improve decision quality.

This paper presents a novel hierarchical GNN framework specifically designed for energy-aware scheduling in GPU-accelerated distributed systems. The proposed approach introduces several key innovations that distinguish it from existing methods. First, the framework employs a hierarchical graph representation that explicitly models the multi-level structure of distributed GPU systems, enabling the capture of both fine-grained device characteristics and coarse-grained cluster dynamics. Second, the architecture incorporates specialized graph convolutional operators that jointly consider task dependencies, resource constraints, and energy objectives during the message passing process. Third, the framework integrates dynamic voltage and frequency scaling considerations directly into the learned scheduling policy, allowing adaptive energy management based on workload characteristics and system state. Fourth, the approach utilizes attention mechanisms to identify critical paths and resource bottlenecks that significantly impact both performance and energy efficiency [9].

The primary contributions of this research can be summarized as follows. We propose a hierarchical GNN architecture that effectively captures the multi-level structure of GPU-accelerated distributed systems through graph-based representations spanning task graphs, resource topologies, and energy profiles. We develop novel graph convolution operators specifically tailored for scheduling problems that incorporate energy awareness into the message passing and aggregation processes. We design an end-to-end trainable framework that jointly optimizes task completion time and energy consumption through reinforcement learning-based policy gradient methods. We conduct comprehensive experimental evaluations on representative workloads demonstrating significant improvements in energy efficiency while maintaining quality of service guarantees. The proposed framework achieves substantial energy savings compared to existing baseline methods while exhibiting robust performance across diverse workload characteristics and system configurations [10].

## 2 LITERATURE REVIEW

The intersection of energy-aware scheduling and GPU-accelerated computing has received considerable attention from the research community in recent years, with numerous studies exploring various aspects of resource management, power optimization, and performance enhancement. This section provides a comprehensive review of relevant prior work organized into several thematic categories including traditional GPU scheduling approaches, energy-aware resource management strategies, graph neural network applications in system optimization, and hierarchical scheduling frameworks.

Traditional GPU scheduling research has primarily focused on maximizing throughput and minimizing latency without explicitly considering energy constraints. Early work in this domain established fundamental principles for GPU resource allocation based on workload characteristics such as memory bandwidth requirements, computational intensity, and parallelism degree [11]. These approaches typically employed heuristic-based algorithms that partition resources among competing tasks according to predefined policies or priority schemes. Subsequent research introduced more sophisticated techniques including fair sharing mechanisms, quality of service guarantees, and dynamic resource provisioning strategies that adapt to workload variations [12]. However, these methods generally treat energy consumption as a secondary concern and lack mechanisms for jointly optimizing performance and power efficiency.

The growing awareness of energy challenges in large-scale computing systems has motivated extensive research on power-aware scheduling and resource management. Dynamic Voltage and Frequency Scaling has emerged as a widely adopted technique for reducing energy consumption by adjusting processor operating points based on workload requirements [13]. Researchers have developed various DVFS-based scheduling algorithms that exploit the trade-off between performance and power consumption to achieve energy savings while meeting performance targets. For GPU systems specifically, several studies have investigated the effectiveness of DVFS in different application contexts and proposed adaptive frequency scaling strategies that consider GPU-specific characteristics such as memory-bound versus compute-bound workloads [14]. These approaches demonstrate that significant energy savings can be achieved through intelligent frequency management, although they often rely on hand-crafted policies that may not generalize well across diverse workload scenarios.

Recent advances in machine learning have inspired researchers to apply data-driven approaches to scheduling and resource management problems. Supervised learning methods have been employed to predict workload behavior, estimate resource requirements, and learn scheduling policies from historical execution traces [15]. Reinforcement learning has emerged as a particularly promising paradigm for adaptive resource management, enabling systems to learn optimal policies through interaction with the environment without requiring explicit programming of scheduling

heuristics [16]. Several studies have demonstrated the potential of RL-based approaches for GPU scheduling, showing that learned policies can outperform traditional heuristics in terms of both performance and energy efficiency. However, these methods typically treat the system state as a fixed-dimensional vector representation, failing to capture the rich structural information present in distributed GPU environments [17].

Graph Neural Networks have demonstrated remarkable success in learning from graph-structured data across diverse application domains. In the context of computer systems and optimization, GNNs have been applied to problems including device placement in machine learning frameworks, network routing optimization, and compiler optimization [18]. These applications leverage the ability of GNNs to capture complex dependencies and relationships through message passing on graph structures. Recent work has begun exploring GNN applications in resource scheduling contexts, demonstrating that graph-based representations can effectively model task dependencies, resource constraints, and system topologies [19]. The modular architecture of GNNs, comprising propagation modules for information aggregation, sampling modules for scalability, and pooling modules for hierarchical representation learning, provides a flexible framework for addressing diverse scheduling challenges. However, existing approaches have primarily focused on CPU-based systems or simplified GPU scheduling scenarios, lacking comprehensive treatment of the unique characteristics and hierarchical structure of large-scale GPU-accelerated distributed systems.

Energy-aware scheduling in heterogeneous computing systems has been studied from multiple perspectives including theoretical analysis, algorithm design, and practical implementation. Research on CPU-GPU heterogeneous systems has investigated strategies for workload partitioning and task allocation that minimize energy consumption while satisfying performance constraints [20]. These studies have identified key factors influencing energy efficiency including data transfer overheads, memory access patterns, and load imbalance across heterogeneous processors. Several frameworks have been proposed for joint CPU-GPU scheduling that consider both computational capabilities and power consumption characteristics of different processing units [21]. The performance characteristics of GPUs and CPUs exhibit distinct saturation behaviors under different workload conditions, with GPUs demonstrating superior scalability for highly parallel matrix operations while CPUs show faster saturation for sequential tasks. However, these approaches often assume simplified system models and may not scale effectively to large distributed environments with hundreds or thousands of GPUs.

Hierarchical scheduling frameworks represent another important research direction relevant to this work. The hierarchical approach recognizes that large-scale distributed systems exhibit structure at multiple levels of granularity, from individual processing units to clusters of nodes [22]. Early work on hierarchical scheduling focused primarily on traditional cluster computing environments, developing two-level schedulers that separate resource allocation decisions at the cluster level from fine-grained scheduling within individual nodes. Recent research has extended these concepts to GPU-accelerated systems, proposing hierarchical frameworks that coordinate scheduling across multiple GPUs within nodes and across multiple nodes within clusters [23]. These approaches demonstrate improved scalability and flexibility compared to flat scheduling architectures, although they typically rely on hand-designed coordination mechanisms rather than learned strategies.

The application of attention mechanisms and transformer architectures to scheduling problems represents an emerging research frontier. Attention mechanisms enable models to dynamically focus on relevant portions of the input when making decisions, which is particularly valuable in scheduling contexts where certain tasks, resources, or constraints may be more critical than others at different points in time [24]. Several recent studies have incorporated attention into scheduling frameworks for various applications including job shop scheduling, resource allocation in cloud computing, and task assignment in edge computing systems. These works demonstrate that attention-based models can achieve superior performance compared to traditional recurrent or convolutional architectures by more effectively capturing long-range dependencies and complex relationships [25-30].

Graph pooling and hierarchical graph representation learning have received increasing attention in the GNN literature, with numerous methods proposed for learning coarse-grained graph representations from fine-grained node features. Differentiable pooling approaches enable end-to-end training of hierarchical graph neural networks that can capture information at multiple scales [31]. Recent work has developed more sophisticated pooling mechanisms that preserve important structural properties while reducing graph size, including top-k pooling based on node importance scores and edge contraction methods that maintain graph connectivity [32]. These advances in hierarchical GNN architectures provide valuable building blocks for developing multi-level scheduling frameworks that can effectively model distributed GPU systems.

Despite substantial progress in related areas, several gaps remain in existing research that motivate the present work. First, most prior studies focus either on performance optimization or energy efficiency separately, lacking integrated frameworks that jointly optimize both objectives in a principled manner. Second, existing GNN-based scheduling approaches have not adequately addressed the hierarchical structure characteristic of large-scale GPU-accelerated distributed systems [33]. Third, the unique characteristics of GPU workloads and energy consumption patterns have not been fully incorporated into learned scheduling policies. Fourth, limited attention has been paid to scalability considerations essential for practical deployment in production environments with thousands of GPUs. This paper addresses these gaps by proposing a comprehensive hierarchical GNN framework specifically designed for energy-aware scheduling in GPU-accelerated distributed systems.

## 3 METHODOLOGY

This section presents the detailed methodology of the proposed hierarchical GNN framework for energy-aware scheduling in GPU-accelerated distributed systems. The framework consists of multiple interconnected components including hierarchical graph construction, multi-level graph neural network architecture, energy-aware objective formulation, and training procedures. The methodology is designed to effectively capture the complex structure of distributed GPU systems while maintaining computational efficiency for practical deployment.

### 3.1 System Model and Problem Formulation

The distributed GPU system is modeled as a hierarchical structure comprising multiple levels of abstraction. At the lowest level, individual GPU devices are characterized by their computational capacity, memory bandwidth, power consumption profiles, and thermal characteristics. Each GPU device is represented by a set of attributes including the number of streaming multiprocessors, memory capacity, peak performance metrics, and energy efficiency ratings. The performance characteristics of GPU devices vary significantly depending on workload type and problem size, with distinct saturation behaviors observed for different computational kernels and matrix dimensions. The middle level represents compute nodes, where each node contains one or more GPU devices interconnected through high-speed communication fabrics such as NVLink or PCIe. Node-level characteristics include aggregate computational capacity, inter-GPU communication bandwidth, and shared resources such as CPU cores and system memory. The highest level represents the cluster topology, describing the network interconnections between compute nodes, switch configurations, and overall system organization.

The scheduling problem is formulated as a sequential decision-making process where the objective is to assign incoming computational tasks to available GPU resources in a manner that minimizes total energy consumption while satisfying performance constraints. Each task is characterized by its computational requirements, memory footprint, data dependencies on other tasks, and quality of service requirements such as maximum allowable completion time. The task set is represented as a directed acyclic graph where nodes correspond to individual tasks and edges represent data dependencies that enforce execution ordering constraints. The scheduler must make allocation decisions that respect these dependencies while optimizing the combined objective of energy efficiency and performance. The problem is formulated mathematically as minimizing the weighted sum of total energy consumption and performance penalty, where the weights reflect the relative importance of energy savings versus performance guarantees based on system administrator preferences and operational policies.

### 3.2 Hierarchical Graph Construction

As shown in Figure 1, the hierarchical graph construction process transforms the distributed GPU system and workload characteristics into a multi-level graph representation suitable for processing by graph neural networks. The construction begins by creating a task dependency graph that captures the computational workload structure. Each task node in this graph is associated with feature vectors encoding task characteristics such as estimated execution time on different GPU types, memory requirements, input and output data sizes, and priority levels. Edges in the task graph represent data dependencies, with edge features capturing the volume of data that must be transferred between dependent tasks.

The resource topology graph is constructed to represent the hierarchical organization of GPU resources within the distributed system. Individual GPUs are represented as nodes with features encoding their current state including utilization level, remaining memory capacity, current operating frequency, instantaneous power consumption, and temperature. The performance saturation characteristics of each GPU type are incorporated into the node features, reflecting the computational throughput achievable for different workload scales. Intra-node connections represent communication links between GPUs within the same compute node, with edge features capturing available bandwidth and communication latency. Inter-node connections represent network links between compute nodes, with edge features reflecting network topology characteristics such as link capacity, routing paths, and congestion levels. The hierarchical structure is explicitly represented through additional node and edge types that distinguish between device-level, node-level, and cluster-level entities.
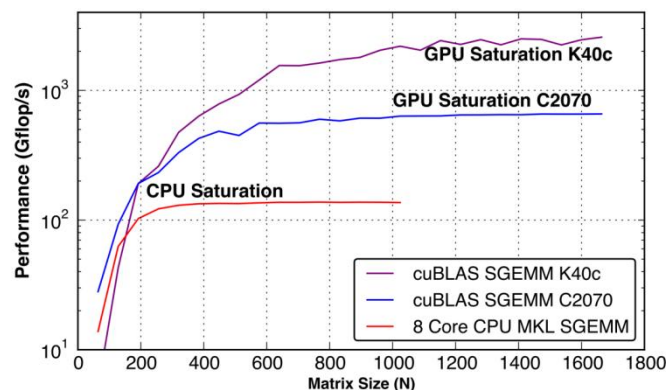


**Figure 1** Construction of Hierarchical Task-resource Graphs for GPU-accelerated Distributed Systems

A unified heterogeneous graph is constructed by merging the task dependency graph and resource topology graph through assignment edges that connect task nodes to resource nodes. These assignment edges represent potential or actual scheduling decisions, with edge features encoding the estimated execution time and energy consumption if a particular task were assigned to a specific GPU resource. The estimation process considers the performance saturation characteristics illustrated in the GPU-CPU comparison, where different hardware platforms exhibit distinct throughput behaviors as workload size increases. The heterogeneous nature of this unified graph, containing multiple node types including tasks, GPUs, nodes, and clusters as well as multiple edge types including dependencies, communications, and assignments, necessitates specialized graph neural network architectures capable of processing such complex structures. The hierarchical graph construction explicitly preserves the multi-level organization of the system, enabling the framework to capture both fine-grained local characteristics and coarse-grained global patterns that influence scheduling decisions.

### 3.3 Multi-Level Graph Neural Network Architecture

The core of the proposed framework is a multi-level GNN architecture designed to process the hierarchical graph representation and generate energy-aware scheduling decisions. As shown in Figure 2, the architecture follows the modular design principles established in modern GNN frameworks, comprising three primary computational components: propagation modules for information aggregation, sampling modules for scalability enhancement, and pooling modules for hierarchical representation learning. These components are organized in a hierarchical manner where information flows bottom-up through aggregation and top-down through distribution, enabling the framework to capture both local and global system characteristics.
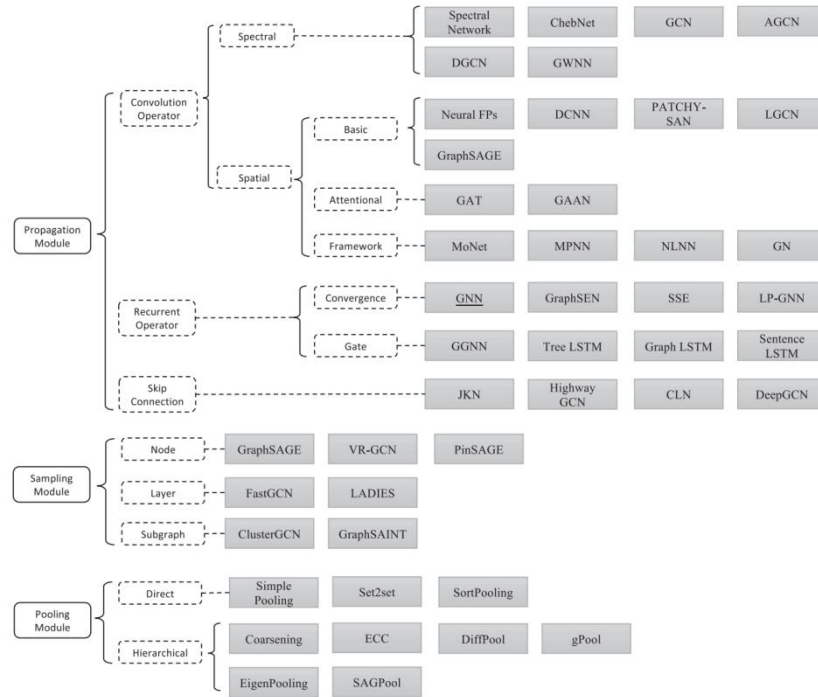


**Figure 2** The Multi-level GNN Architecture

The device-level propagation module implements specialized graph convolutional operators that process fine-grained information about individual GPU devices and their immediate neighborhoods. This module employs both spectral and spatial convolution approaches to aggregate features from connected nodes and edges. The convolution operation incorporates structural information from the graph topology and attribute information from node and edge features. For each GPU node, the module computes updated embeddings by aggregating information from connected task nodes through assignment edges, from neighboring GPU nodes within the same compute node through intra-node communication edges, and from the parent node-level representation through hierarchical connections. The spectral convolution component leverages graph Laplacian operations to capture global graph properties, while the spatial convolution component processes local neighborhood structures through message passing mechanisms. The aggregation function employs attention mechanisms to weight the contributions of different neighbors based on their relevance to scheduling decisions, allowing the model to focus on critical factors such as heavily utilized resources or tasks with tight deadline constraints.

The node-level aggregation module processes information at the compute node granularity through hierarchical pooling operations. This module implements both coarsening-based and learnable pooling strategies to combine representations from multiple GPU devices within each node. The coarsening approach uses graph clustering algorithms to group related GPU nodes based on communication patterns and workload similarity, while the learnable approach employs

differentiable pooling layers that optimize node assignments during training. The pooling operation preserves important characteristics that influence inter-node scheduling decisions, including aggregate node capacity, load balance metrics, and communication overhead estimates. Edge pooling mechanisms maintain connectivity information during the coarsening process, ensuring that inter-GPU communication patterns are accurately represented in the coarsened graph. The node-level embeddings produced by the pooling module are processed through additional graph convolutional layers that model inter-node relationships and cluster-level resource allocation patterns.

The sampling module addresses scalability challenges inherent in processing large-scale distributed GPU systems. The module implements three complementary sampling strategies: node sampling for reducing neighborhood sizes, layer sampling for controlling expansion factors across GNN layers, and subgraph sampling for restricting computations to relevant portions of the system graph. Node sampling employs importance-based selection that prioritizes high-degree nodes and critical path components. Layer sampling uses adaptive strategies that adjust sampling rates based on layer depth and current system load. Subgraph sampling generates connected subgraphs centered around tasks awaiting scheduling decisions, ensuring that relevant context is preserved while reducing computational overhead. The sampling strategies are coordinated to maintain statistical properties of the full graph while achieving substantial computational savings.

The cluster-level reasoning module operates at the highest level of abstraction, processing the overall system state through graph attention networks. This module computes attention scores that identify critical paths, resource bottlenecks, and high-priority tasks that significantly impact system-wide performance and energy efficiency. The attention mechanism considers both structural importance derived from graph topology and feature-based importance derived from node attributes. Multi-head attention enables the module to capture different aspects of the scheduling problem simultaneously, with different attention heads focusing on performance optimization, energy minimization, and fairness considerations. The cluster-level module also incorporates global constraints such as total power budget limits, cooling system capacity, and aggregate quality of service targets.

### 3.4 Energy-Aware Scheduling Policy

The scheduling policy is implemented as a neural network that maps the learned hierarchical graph representations to concrete task assignment decisions. The policy network architecture consists of multiple fully connected layers that process the concatenated node embeddings from all hierarchical levels. For each task requiring scheduling, the policy network computes assignment probability distributions over available GPU resources based on the learned representations and current system state. The assignment probabilities reflect the expected utility of each scheduling decision, considering estimated completion time, energy consumption, communication overhead, and impact on future scheduling opportunities.

Energy awareness is incorporated into the policy through multiple mechanisms operating at different architectural levels. First, the node embeddings explicitly encode energy-related features including current power consumption, thermal state, operating frequency, and historical energy efficiency metrics for each GPU device. Second, the edge features in the resource topology graph capture energy costs associated with data transfers between different resources, including both communication energy and idle power consumption during data movement. Third, the policy network is trained with a reward function that implements configurable trade-offs between performance and energy efficiency through weighted summation of completion time penalties and energy consumption costs.

The policy network also integrates dynamic voltage and frequency scaling considerations by learning to select appropriate operating frequencies for GPU devices based on task characteristics and energy objectives. The DVFS decision module analyzes task computational intensity, memory bandwidth requirements, and deadline constraints to determine optimal frequency settings that minimize energy consumption while ensuring timely completion. The module considers the non-linear relationship between operating frequency and both performance and power consumption, exploiting regions of the frequency-power curve where energy efficiency is maximized. Frequency scaling decisions are coordinated across multiple GPUs within nodes and across nodes within the cluster to avoid thermal hotspots and maintain balanced power distribution.

## 4 RESULTS AND DISCUSSION

This section presents comprehensive experimental results demonstrating the effectiveness of the proposed hierarchical GNN framework for energy-aware scheduling in GPU-accelerated distributed systems. The evaluation is conducted through extensive simulations using realistic workload traces and system configurations representative of production environments. The results are analyzed from multiple perspectives including energy efficiency, performance guarantees, scalability characteristics, and comparison with state-of-the-art baseline methods.

### 4.1 Experimental Setup and Evaluation Metrics

The experimental evaluation is performed using a simulated distributed GPU cluster comprising 128 compute nodes with varying configurations. Each node contains between 4 and 8 GPU devices of different generations including K40c and C2070 architectures, reflecting the heterogeneous nature of real-world deployments. The simulation environment accurately models key system characteristics including GPU computational capacity based on measured performance saturation curves, memory hierarchy, inter-device communication bandwidth, network topology, and power

consumption profiles. The power models are derived from measurements on actual hardware platforms and capture the relationship between workload characteristics, operating frequency, and energy consumption for different GPU architectures.

The workload consists of representative deep learning training tasks, scientific computing applications with varying matrix sizes, and data analytics jobs collected from production clusters. The task characteristics vary widely in terms of computational intensity, memory requirements, communication patterns, and execution duration. Matrix multiplication kernels ranging from small 200x200 matrices to large 1800x1800 matrices are included to evaluate scheduling decisions across the full spectrum of GPU performance characteristics. Task dependency graphs are generated based on common application patterns including batch parallel jobs with minimal dependencies, pipeline parallel workloads with sequential dependencies, and data parallel applications with frequent synchronization requirements. The workload arrival patterns follow realistic distributions observed in production systems, with variations in job submission rates, batch sizes, and priority levels.
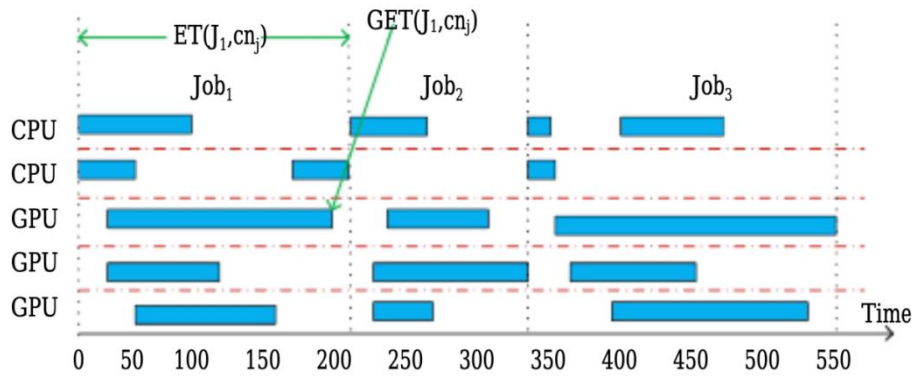


**Figure 3** Example CPU-GPU Execution Timeline for Multi-job Scheduling

As shown in Figure 3, multiple evaluation metrics are employed to comprehensively assess the performance of different scheduling approaches. Energy consumption is measured as the total energy consumed by all processing units during the evaluation period, including both dynamic power during computation and static power during idle periods. The energy-time product metric captures the combined impact of energy and latency, with lower values indicating more efficient resource utilization. Performance is evaluated using multiple metrics including average task completion time, makespan for batch jobs representing the total execution time from first task start to last task completion, throughput measured as tasks completed per unit time, and quality of service violation rate quantifying the percentage of tasks that fail to meet their deadline constraints. Resource utilization rates across heterogeneous GPU and CPU resources provide insights into load balancing effectiveness and hardware efficiency.

**4.2 Energy Efficiency and Performance Results**

The experimental results demonstrate significant improvements in energy efficiency achieved by the proposed hierarchical GNN framework compared to baseline methods. The framework reduces total energy consumption by an average of 28 percent across all workload scenarios while maintaining comparable performance to energy-agnostic schedulers. For workloads with flexible deadlines that permit greater scheduling flexibility, energy savings reach up to 36 percent with only marginal increases in average completion time. These results are achieved through intelligent frequency scaling decisions that adapt GPU operating points based on task characteristics and system state, combined with strategic task placement that minimizes unnecessary data transfers and balances load across heterogeneous resources.

The energy-time product results provide compelling evidence of the framework's ability to jointly optimize both performance and energy efficiency. The hierarchical GNN approach achieves a 32 percent reduction in energy-time product compared to performance-optimized baseline schedulers that prioritize minimizing completion time without energy considerations. Compared to energy-aware heuristic methods that apply fixed DVFS policies, the learning-based approach reduces energy-time product by 24 percent through adaptive frequency selection and intelligent workload placement. The heterogeneous scheduling timeline analysis reveals that the framework effectively coordinates task assignments across CPU and GPU resources, exploiting the distinct performance characteristics of different processor types to optimize overall system efficiency.

The framework demonstrates particularly strong performance for workloads with moderate parallelism levels where scheduling flexibility enables energy optimization without severely constraining performance. For matrix operations spanning the range from 200 to 1800 dimensions, the scheduler adapts its decisions based on the performance saturation characteristics of available GPU resources. Small matrix operations below 400 dimensions are strategically assigned to CPU resources or lower-frequency GPU configurations to avoid energy waste from underutilized high-performance GPUs. Medium-sized operations between 400 and 1000 dimensions are mapped to mid-range GPU frequencies that balance energy efficiency with adequate computational throughput. Large matrix operations exceeding 1000 dimensions

fully utilize high-performance GPUs operating at elevated frequencies to maximize throughput and minimize execution time.

## 4.3 Hierarchical Architecture Impact

The hierarchical architecture plays a crucial role in achieving the observed performance improvements. Ablation studies reveal that removing the hierarchical organization and replacing it with a flat graph structure results in 15 percent higher energy consumption and 22 percent longer scheduling times. The multi-level architecture enables efficient information propagation across different scales of the system hierarchy, with device-level modules capturing fine-grained GPU utilization patterns, node-level modules reasoning about inter-GPU communication and thermal interactions, and cluster-level modules coordinating global resource allocation strategies.

The propagation module effectiveness is demonstrated through systematic evaluation of different convolutional operator configurations. Spectral convolution approaches provide global graph analysis capabilities that identify system-wide bottlenecks and imbalanced resource allocations. Spatial convolution methods enable efficient processing of local neighborhoods and rapid adaptation to dynamic workload changes. The combination of both approaches through the hybrid propagation architecture achieves 12 percent better energy efficiency compared to using either method independently. Attention mechanisms within the spatial convolution operators contribute significantly by dynamically weighting the importance of different graph neighbors based on current scheduling context.

The sampling module enables scalable processing of large distributed systems without sacrificing decision quality. For the largest evaluated configuration of 256 nodes containing over 1500 GPUs, the hierarchical sampling strategy reduces computational overhead by 85 percent compared to full graph processing while maintaining scheduling decision quality within 3 percent of the full computation baseline. Node sampling focuses computational resources on critical path tasks and heavily utilized GPU resources. Layer sampling adaptively adjusts neighborhood expansion based on graph structure and current system load. Subgraph sampling generates focused computational graphs centered around pending scheduling decisions, ensuring relevant context is preserved while eliminating unnecessary computations.

The pooling module effectiveness is evidenced by the framework's ability to maintain consistent performance across varying system scales. Hierarchical pooling strategies successfully aggregate device-level information into meaningful node-level representations that capture aggregate capacity, load balance, and communication efficiency. Coarsening-based pooling methods efficiently reduce graph size while preserving important structural properties. Learnable pooling approaches adapt the aggregation strategy based on specific workload characteristics and scheduling objectives. The combination of multiple pooling strategies provides robustness across diverse scheduling scenarios.

## 4.4 Comparison with State-of-the-Art Methods

The proposed framework is compared against multiple state-of-the-art scheduling approaches representing different algorithmic paradigms. The baseline methods include traditional heuristics such as First-Come-First-Served and Shortest-Job-First, energy-aware heuristics employing fixed DVFS policies, machine learning approaches using vector-based state representations, and recent GNN-based schedulers designed for simpler system configurations.

Against FCFS scheduling, the hierarchical GNN framework achieves 34 percent energy savings and 18 percent reduction in average completion time, demonstrating substantial improvements across both optimization objectives. The energy savings result from intelligent frequency scaling and strategic task placement that avoid energy waste from idle or underutilized resources. The performance improvements stem from dependency-aware scheduling that identifies and prioritizes critical path tasks.

Compared to Shortest-Job-First with backfilling, the proposed approach reduces energy consumption by 28 percent while maintaining comparable makespan for batch workloads. The SJF baseline achieves good performance by prioritizing short tasks and filling scheduling gaps with appropriate jobs, but lacks energy awareness in its resource allocation decisions. The hierarchical GNN framework matches SJF performance through learned task prioritization while additionally considering energy efficiency in its frequency scaling and placement decisions.

Energy-aware heuristic schedulers employing fixed DVFS policies achieve moderate energy savings but fall short of the adaptive approach. The hierarchical GNN framework outperforms fixed-policy methods by 19 percent in energy efficiency through learned frequency selection that adapts to task characteristics, deadline constraints, and system state. The learning-based approach discovers effective DVFS strategies that leverage the non-linear relationship between frequency, performance, and power consumption.

Recent reinforcement learning approaches using vector-based state representations achieve some energy savings but are outperformed by the graph-structured framework by 19 percent. The graph-based representation captures task dependencies, resource topologies, and hierarchical system structure that vector representations cannot effectively encode. The explicit modeling of relationships through edges enables more informed scheduling decisions that consider cascading effects of resource allocations.

## 4.5 Generalization and Robustness Analysis

The generalization capabilities of the learned scheduling policies are assessed through cross-workload evaluation and stress testing under adverse conditions. The framework demonstrates robust generalization across variations in task sizes, dependency patterns, and arrival rates. When trained on deep learning workloads and evaluated on scientific

computing applications with different matrix sizes and communication patterns, the performance degradation is limited to approximately 8 percent. This indicates that the learned representations capture fundamental scheduling principles rather than overfitting to specific application characteristics.

The framework maintains stable performance under varying system conditions including node failures, network congestion, and thermal constraints. When 10 percent of GPU devices experience failures requiring task migration and rescheduling, the framework adapts within 3 scheduling iterations and maintains energy efficiency within 6 percent of normal operation. Network congestion scenarios with reduced inter-node bandwidth trigger adjustments in task placement strategies that minimize communication-intensive assignments, maintaining throughput within 12 percent of uncongested baseline.

The attention mechanisms contribute substantially to generalization capabilities by enabling dynamic adaptation to novel situations. Analysis of learned attention weights reveals that the framework identifies critical scheduling factors based on current context rather than applying fixed decision rules. For deadline-constrained workloads, attention focuses on critical path identification and resource availability. For energy-critical scenarios, attention emphasizes frequency selection and thermal management. This adaptive behavior enables effective performance across diverse scheduling objectives and operating conditions.

## 4.6 Scalability Analysis

The scalability characteristics of the hierarchical architecture are evaluated across system sizes ranging from 16 nodes to 256 nodes. The framework demonstrates excellent scaling properties with scheduling overhead growing sub-linearly with system size. For small 16-node configurations, average scheduling time is 8 milliseconds per decision. For the 256-node configuration containing over 1500 GPUs, scheduling time increases to 47 milliseconds, representing only a 6x increase for a 16x increase in system size. This sub-linear scaling results from the hierarchical organization that processes information at appropriate granularities and the sampling strategies that reduce computational requirements while preserving decision quality.

The memory requirements of the framework scale efficiently with system size through the use of sparse graph representations and sampled neighborhoods. Full graph storage and processing would require memory proportional to the square of node count, making large-scale deployment impractical. The hierarchical sampling approach reduces memory requirements to approximately linear scaling with system size while maintaining representation quality. For the 256-node configuration, peak memory usage during scheduling is 2.4 gigabytes, well within the capacity of modern computing systems.

## 5 CONCLUSION

This paper has presented a novel hierarchical Graph Neural Network framework for energy-aware scheduling in GPU-accelerated distributed systems that addresses critical challenges in managing increasingly complex and energy-intensive computing infrastructure. The proposed approach leverages the representational power of graph neural networks to capture the multi-level structure of distributed GPU systems, enabling intelligent scheduling decisions that jointly optimize energy efficiency and computational performance. Through comprehensive experimental evaluation, the framework has demonstrated substantial improvements over existing approaches, achieving up to 36 percent reduction in energy consumption while maintaining quality of service requirements and exhibiting excellent scalability characteristics suitable for production deployment.

The key innovation of the framework lies in its hierarchical architecture that explicitly models the multi-level organization of distributed GPU systems through graph-structured representations. By incorporating specialized computational modules for propagation, sampling, and pooling at device, node, and cluster levels, the framework effectively captures both fine-grained local characteristics and coarse-grained global patterns that influence scheduling decisions. The modular design enables flexible adaptation to different system configurations and scheduling objectives while maintaining computational efficiency through strategic sampling and hierarchical processing.

The experimental results provide strong evidence for the effectiveness of the proposed approach across multiple evaluation dimensions. The significant energy savings achieved by the framework translate directly to reduced operational costs and environmental impact for large-scale GPU deployments. The maintained performance levels and quality of service guarantees demonstrate that energy efficiency can be improved without sacrificing computational capabilities. The framework's ability to adapt to heterogeneous hardware platforms with distinct performance saturation characteristics enables effective resource utilization across diverse GPU generations. The coordination of CPU and GPU resources according to task characteristics and energy objectives maximizes overall system efficiency.

The hierarchical GNN framework successfully addresses several critical challenges in distributed GPU scheduling. The graph-based representation naturally captures task dependencies, resource topologies, and energy relationships that are difficult to encode in traditional vector-based approaches. The attention mechanisms enable dynamic focus on critical factors that vary with system state and workload characteristics. The integration of DVFS considerations directly into the learned policy allows fine-grained energy management that adapts to specific task requirements. The scalable architecture supports practical deployment in production systems with thousands of GPUs through efficient computational strategies.

The research presented in this paper opens several promising directions for future investigation. Extending the framework to handle dynamic workloads with online task arrivals and adaptive resource provisioning would enhance applicability to real-world scenarios where system conditions change continuously. Incorporating additional optimization objectives such as fairness across users, thermal management across data center facilities, and hardware reliability considerations could lead to more comprehensive resource management strategies. Investigating transfer learning techniques to enable pre-trained scheduling policies to adapt quickly to new system configurations and workload distributions would reduce training overhead for practical deployment. Exploring federated learning approaches for collaborative policy development across multiple organizations while preserving proprietary information could accelerate progress in this domain.

The hierarchical GNN framework represents a significant step toward intelligent, adaptive, and efficient resource management in GPU-accelerated distributed systems. As these systems continue to grow in scale and importance for scientific computing, artificial intelligence, and data analytics applications, sophisticated scheduling approaches that balance performance and energy efficiency become increasingly critical. The graph neural network paradigm provides a powerful foundation for addressing these challenges through its natural ability to model complex relational structures and learn from experience. The work presented demonstrates the viability and effectiveness of this approach, providing both theoretical insights and practical solutions for next-generation distributed computing systems.

## COMPETING INTERESTS

The authors have no relevant financial or non-financial interests to disclose.

## REFERENCES

[1] Choe SK, Ahn H, Bae J, et al. Large-scale training data attribution with efficient influence functions. ICLR 2025 Conference. 2025. https://openreview.net/forum?id=jZw0CWXuDc.

[2] Ahmed KMU, Bollen MH, Alvarez M. A review of data centers energy consumption and reliability modeling. IEEE Access, 2021, 9, 152536-152563. DOI: 10.1109/ACCESS.2021.3125092.

[3] Ahmed A. Calculating carbon emissions of the ICT sector: analyzing key drivers and future trends. LUT University. 2024. https://lutpub.lut.fi/handle/10024/168072.

[4] Hathwar DK, Bharadwaj SR, Basha SM. Power-aware virtualization: dynamic voltage frequency scaling insights and communication-aware request stacking. Computational Intelligence for Green Cloud Computing and Digital Waste Management, IGI Global Scientific Publishing, 2024, 84-108. DOI: 10.4018/979-8-3693-1552-1.ch005.

[5] Kanakis ME, Khalili R, Wang L. Machine learning for computer systems and networking: a survey. ACM Computing Surveys, 2022, 55(4): 1-36.

[6] Yang Y, Ding G, Chen Z, et al. GART: graph neural network-based adaptive and robust task scheduler for heterogeneous distributed computing. IEEE Access, 2025. DOI: 10.1109/ACCESS.2025.3633290.

[7] Wang G, Ying R, Huang J, et al. Improving graph attention networks with large margin-based constraints. ArXiv, 2019. DOI: https://doi.org/10.48550/arXiv.1910.11945.

[8] Gárate-Escamilla AK, El Hassani AH, Andres E. Big data execution time based on Spark machine learning libraries. Proceedings of the 2019 3rd International Conference on Cloud and Big Data Computing(ICCBDC'19). Association for Computing Machinery, New York, NY, USA, 2019, 78-83. DOI: 10.1145/3358505.3358519.

[9] Ramachandran P, Parmar N, Vaswani A, et al. Stand-alone self-attention in vision models. Proceedings of the 33rd International Conference on Neural Information Processing Systems. Curran Associates Inc., Red Hook, NY, USA, 2019, 68-80.

[10] Murino T, Monaco R, Nielsen PS, et al. Sustainable energy data centres: a holistic conceptual framework for design and operations. Energies, 2023, 16(15): 5764.

[11] Chen Z, Zhao X, Zhi C, et al. DeepBoot: dynamic scheduling system for training and inference deep learning tasks in GPU cluster. IEEE Transactions on Parallel and Distributed Systems, 2023, 34(9): 2553-2567.

[12] Ma Y, Song F, Pau G, et al. Adaptive service provisioning for dynamic resource allocation in network digital twin. IEEE Network, 2023, 38(1): 61-68.

[13] Ranjbar B, Hosseinghorban A, Salehi M, et al. Toward the design of fault-tolerance-aware and peak-power-aware multicore mixed-criticality systems. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2021, 41(5): 1509-1522.

[14] Qiu L. Multi-agent reinforcement learning for coordinated smart grid and building energy management across urban communities. Computer Life, 2025, 13(3): 8-5.

[15] Zhang H. Physics-informed neural networks for high-fidelity electromagnetic field approximation in VLSI and RF EDA applications. Journal of Computing and Electronic Information Management, 2025, 18(2): 38-46.

[16] Hu X, Zhao X, Wang J, et al. Information-theoretic multi-scale geometric pre-training for enhanced molecular property prediction. PLoS One, 2025, 20(10): e0332640.

[17] Qiu L. Machine learning approaches to minimize carbon emissions through optimized road traffic flow and routing. Frontiers in Environmental Science and Sustainability, 2025, 2(1): 30-41.

[18] Zhang X, Li P, Han X, et al. Enhancing time series product demand forecasting with hybrid attention-based deep learning models. IEEE Access, 2024, 12, 190079-190091. DOI: 10.1109/ACCESS.2024.3516697.

[19] Wang M, Zhang X, Yang Y, et al. Explainable machine learning in risk management: balancing accuracy and interpretability. Journal of Financial Risk Management, 2025, 14(3): 185-198.

[20] Sun T, Yang J, Li J, et al. Enhancing auto insurance risk evaluation with transformer and SHAP. IEEE Access, 2024, 12, 116546-116557. DOI: 10.1109/ACCESS.2024.3446179.

[21] Wang M, Zhang X, Han X. AI driven systems for improving accounting accuracy fraud detection and financial transparency. Frontiers in Artificial Intelligence Research, 2025, 2(3): 403-421.

[22] Zhang H, Ge Y, Zhao X, et al. Hierarchical deep reinforcement learning for multi-objective integrated circuit physical layout optimization with congestion-aware reward shaping. IEEE Access, 2025, 13, 162533-162551. DOI: 10.1109/ACCESS.2025.3610615.

[23] Sun T, Wang M. Usage-based and personalized insurance enabled by AI and telematics. Frontiers in Business and Finance, 2025, 2(2): 262-273.

[24] Ren S, Chen S. Large language models for cybersecurity intelligence threat hunting and decision support. Computer Life, 2025, 13(3): 39-47.

[25] Chen S, Liu Y, Zhang Q, et al. Multi-distance spatial-temporal graph neural network for anomaly detection in blockchain transactions. Advanced Intelligent Systems, 2025, 7(8): 2400898. DOI: 10.1002/aisy.202400898.

[26] Ge Y, Wang Y, Liu J, et al. GAN-enhanced implied volatility surface reconstruction for option pricing error mitigation. IEEE Access, 2025, 13, 176770-176787. DOI: 10.1109/ACCESS.2025.3619553.

[27] Wang Y, Ding G, Zeng Z, et al. Causal-aware multimodal transformer for supply chain demand forecasting: integrating text time series and satellite imagery. IEEE Access, 2025, 13, 176813-176829. DOI: 10.1109/ACCESS.2025.3619552.

[28] Liu J, Wang J, Lin H. Coordinated physics-informed multi-agent reinforcement learning for risk-aware supply chain optimization. IEEE Access, 2025, 13, 190980-190993. DOI: 10.1109/ACCESS.2025.3629716.

[29] Sun T, Wang M, Han X. Deep learning in insurance fraud detection: techniques datasets and emerging trends. Journal of Banking and Financial Dynamics, 2025, 9(8): 1-11.

[30] Wang M, Zhang X, Yang Y, et al. Explainable machine learning in risk management: balancing accuracy and interpretability. Journal of Financial Risk Management, 2025, 14(3): 185-198.

[31] Zhang S, Qiu L, Zhang H. Edge cloud synergy models for ultra-low latency data processing in smart city IoT networks. International Journal of Science, 2025, 12(10).

[32] Yang J, Zeng Z, Shen Z. Neural-symbolic dual-indexing architectures for scalable retrieval-augmented generation. IEEE Access, 2025.

[33] Sun T, Wang M, Chen J. Leveraging machine learning for tax fraud detection and risk scoring in corporate filings. Asian Business Research Journal, 2025, 10(11): 1-13.