

# Enhancing Domain-Specific Language Models with Knowledge Graph Injection and Graph Attention Networks

Jean Dupont, Marie Laurent, Michael Reynolds

Department of Mathematical and Computational Sciences, University of Toronto Mississauga,  
Mississauga ON L5L 1C6, Canada

## Abstract

The rapid evolution of Large Language Models has revolutionized natural language processing, yet these models frequently exhibit limitations when deployed in specialized high-stakes domains such as medicine, law, and engineering. A primary deficiency is the propensity for hallucination and the inability to access up-to-date, structured factual knowledge that was not present or emphasized during the pre-training phase. This paper proposes a novel architecture that integrates Domain-Specific Knowledge Graphs with pre-trained language models utilizing Graph Attention Networks. By employing a dual-stream mechanism that processes textual input alongside structured graph data, we facilitate a deep injection of semantic relationships into the latent space of the language model. The Graph Attention Network component dynamically weighs the importance of neighboring entities within the knowledge graph, allowing the model to attend to the most relevant factual context corresponding to the input query. We evaluate this approach on two distinct domain-specific datasets involving biomedical and legal texts. Our experimental results demonstrate that this injection mechanism significantly outperforms standard fine-tuning approaches in terms of factual accuracy and reasoning capabilities. The proposed method offers a scalable pathway toward creating more reliable and logically sound domain-specific artificial intelligence systems.

## Keywords

Knowledge Graphs, Graph Attention Networks, Large Language Models, Domain Adaptation

## 1 Introduction

The advent of transformer-based architectures has ushered in a new era of capability in artificial intelligence, particularly in the generation and understanding of human language. Models trained on vast corpora of general internet text have demonstrated remarkable proficiency in varied tasks, ranging from translation to creative writing. However, the stochastic nature of these models presents substantial challenges when they are applied to vertical domains requiring high precision and strict adherence to factual reality. In fields such as healthcare, finance, and jurisprudence, the cost of errors is prohibitively high. The phenomenon known as hallucination, where a model generates plausible-sounding but factually incorrect information, remains a critical bottleneck for the deployment of generative artificial intelligence in professional settings. This issue stems largely from the fact that language models store knowledge implicitly within their parameters, making it difficult to verify, update, or reason over explicit relationships between entities [1]. To mitigate these limitations, researchers have increasingly turned to neuro-symbolic approaches that combine the statistical power of neural networks with the structured reliability of symbolic knowledge bases. Knowledge Graphs serve as an ideal repository for such structured information,

representing entities as nodes and their relationships as edges. While Knowledge Graphs provide a rich source of factual grounding, effectively integrating this discrete structure into the continuous vector space of a neural language model is a non-trivial engineering and theoretical challenge. Early attempts often relied on simple concatenation of retrieved facts to the input text, a method that is computationally inefficient and often fails to capture the complex, multi-hop reasoning required for sophisticated queries [2]. This paper introduces a robust framework for Enhancing Domain-Specific Language Models with Knowledge Graph Injection and Graph Attention Networks. Unlike static embedding approaches, our method utilizes Graph Attention Networks to process the local neighborhood of entities identified in the input text. By leveraging the attention mechanism inherent in these networks, our model learns to assign varying levels of importance to different nodes in the graph, thereby filtering out irrelevant noise and focusing on the semantic relationships that are most pertinent to the current context. This graph-encoded context is then fused with the textual representations of the language model through a specialized cross-attention injection layer [3]. The contributions of this study are manifold. First, we provide a formalized architecture for the seamless integration of graph-based features into transformer decoders. Second, we demonstrate that the use of attention mechanisms over graph structures allows for better handling of knowledge heterogeneity and sparsity compared to traditional graph convolutional networks. Third, we present empirical evidence that our approach not only improves accuracy but also enhances the interpretability of the model's decision-making process by highlighting the specific graph entities that influenced the output. This research bridges the gap between unstructured textual learning and structured knowledge representation, paving the way for more trustworthy domain-expert systems.

## 2. Related Work

The trajectory of natural language processing has shifted dramatically from rule-based systems to statistical models and, more recently, to deep learning architectures. The introduction of the Transformer architecture marked a pivotal moment, enabling the parallel processing of sequences and the capture of long-range dependencies. Despite their success, purely data-driven models often struggle with tasks requiring external knowledge that is not frequently represented in the training corpus. This section reviews the historical progression of language modeling, the development of knowledge representation, and the intersection of graph neural networks with textual processing.

### 2.1 Language Models and Domain Adaptation

Pre-trained language models, such as BERT and its successors, have achieved state-of-the-art results on the General Language Understanding Evaluation benchmarks. These models learn contextual representations of words by predicting masked tokens or the next token in a sequence. While highly effective for general tasks, their performance degrades in specialized domains where the vocabulary and semantic structures differ significantly from the general web text used for pre-training. Conventional domain adaptation involves fine-tuning the model on a smaller, domain-specific corpus. While this adjusts the linguistic style of the model, it does not necessarily imbue it with the structured logical constraints governing that domain [4]. Furthermore, fine-tuning on small datasets can lead to catastrophic forgetting, where the model loses its general reasoning capabilities. Recent studies have highlighted that simple fine-tuning is insufficient for ensuring factual consistency, necessitating external augmentation strategies [5].

## 2.2 Knowledge Graph Integration Strategies

Knowledge Graphs offer a structured representation of facts, typically stored as triples consisting of a subject, predicate, and object. Integrating these triples into neural networks has been a subject of extensive research. Early methods utilized knowledge graph embeddings, such as TransE or DistMult, to convert entities and relations into vector representations. These pre-computed embeddings were then concatenated with word embeddings in the input layer of the language model. However, this approach treats knowledge integration as a static process, ignoring the context-dependent nature of information relevance [6]. More advanced techniques have explored the use of retrieval-augmented generation, where relevant documents or graph sub-graphs are retrieved and prepended to the prompt. While effective, retrieval-augmented generation often suffers from context window limitations and the retrieval of irrelevant information that distracts the model. The challenge remains to integrate knowledge in a way that is both deep, affecting the internal reasoning of the model, and dynamic, adjusting to the specific query at hand [7]. Approaches that attempt to linearize graphs into text sequences for model consumption often lose the structural topology that makes graphs unique, leading to sub-optimal reasoning performance [8].

## 2.3 Graph Neural Networks in NLP

Graph Neural Networks have emerged as a powerful tool for processing non-Euclidean data structures. In the context of natural language processing, Graph Neural Networks have been used to model syntactic dependency trees, semantic role labeling, and co-reference resolution. Graph Convolutional Networks operate by aggregating information from a node's immediate neighbors to update its representation. However, Graph Convolutional Networks typically assign equal or statically defined weights to all neighbors, which can be problematic in knowledge graphs where a single node may have hundreds of connections, only a few of which are relevant to a specific context [9]. To address this, Graph Attention Networks were introduced, incorporating an attention mechanism that learns to weigh the contribution of each neighbor dynamically. This allows the model to focus on the most informative parts of the graph. In the domain of question answering, Graph Attention Networks have been used to reason over knowledge bases to derive answers that are not explicitly stated in the text. Despite these advancements, the integration of Graph Attention Network outputs into the deep layers of Generative Pre-trained Transformers remains an under-explored area. Most existing hybrid models use shallow fusion techniques at the input or output level. Our work distinguishes itself by employing a deep injection mechanism where graph-aware representations modulate the self-attention matrices of the language model layers [10].

## 3. Methodology

Our proposed framework is designed to synergize the generative fluency of Large Language Models with the factual rigidity of Domain-Specific Knowledge Graphs. The architecture consists of three primary components: the Knowledge Graph Construction and Retrieval module, the Graph Attention Network encoder, and the Multi-Modal Injection Layer. This section details the theoretical underpinnings and operational mechanics of each component.

Figure 1: Architectural Diagram

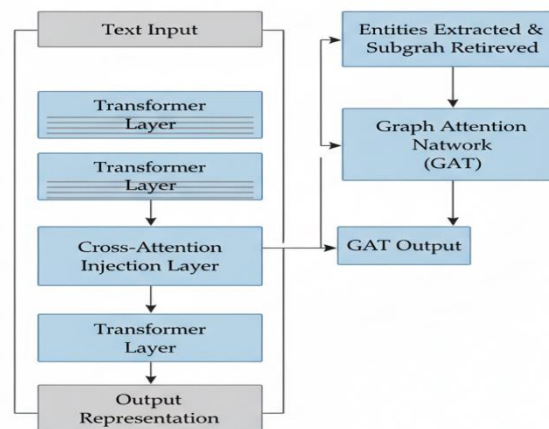


Figure 1: Architectural Diagram

### 3.1 Knowledge Graph Construction and Retrieval

The foundation of our approach is a high-quality domain-specific Knowledge Graph. For the purpose of this study, we utilize existing ontologies relevant to the target domains, specifically the Unified Medical Language System for the biomedical domain and a custom-built legal ontology for the legal domain. The graph is formally defined as a set of vertices and edges, where vertices represent entities and edges represent semantic relations. Given an input text sequence, we first employ an entity linker to identify mentions of graph entities within the text. This process involves Named Entity Recognition followed by disambiguation to map text spans to unique identifiers in the Knowledge Graph. Once the entities are identified, we extract a subgraph centered around these entities. To capture sufficient context without introducing excessive noise, we include all nodes within a 2-hop neighborhood of the identified entities. This subgraph serves as the input to the graph processing module. The limitation to 2-hops is a design choice balancing computational efficiency with the need for multi-step reasoning capabilities [11].

### 3.2 Graph Attention Network Encoder

The extracted subgraph is processed using a Graph Attention Network. Unlike standard Graph Convolutional Networks, the Graph Attention Network computes the hidden states of each node by attending to its neighbors. The core idea is to compute an attention coefficient that indicates the importance of a neighbor node to a central node. This is achieved through a shared linear transformation applied to every node, followed by a self-attention mechanism. The attention mechanism creates a weighted sum of neighbor features. If a node has a high degree of connectivity, the attention mechanism effectively filters out irrelevant connections based on the current feature states. We employ a multi-head attention structure within the graph network to stabilize the learning process. Each head computes independent attention coefficients, and their outputs are concatenated to form the final node representation. This allows the model to capture different types of relationships simultaneously—for example, one head might focus on causal relationships while another focuses on hierarchical classifications. The output of this stage is a set of context-enriched node embeddings that encapsulate both the intrinsic properties of the entities and their local structural context [12].

### 3.3 Deep Knowledge Injection Mechanism

The integration of the graph embeddings into the language model is achieved through a novel injection layer inserted between the transformer blocks of the language model. Standard language models utilize self-attention where tokens attend to other tokens in the sequence. We introduce a cross-attention layer where the query vectors are derived from the text embeddings, and the key and value vectors are derived from the graph node embeddings produced by the Graph Attention Network. This configuration allows the text tokens to query the knowledge graph for relevant information. For instance, if the text mentions a specific drug, the attention mechanism will likely assign high weights to the graph nodes representing that drug's side effects or contraindications, provided they exist in the extracted subgraph. This retrieved graph information is then added to the text representation via a residual connection and layer normalization. By interleaving these injection layers at multiple depths of the language model, we ensure that the generated text is consistently guided by the structured knowledge throughout the generation process. This deep fusion strategy is superior to input-level concatenation because it allows the model to refine its understanding of the graph context as the abstraction level of the text representation increases through the layers [13].

## 4. Experimental Setup

To validate the efficacy of our proposed architecture, we conducted rigorous experiments comparing our model against several strong baselines. The experiments were designed to assess both the factual accuracy of the generated text and the model's ability to utilize the injected knowledge for reasoning tasks.

### 4.1 Datasets and Baselines

We utilized two primary datasets representing distinct domains. The first is MedQA, a large-scale biomedical question-answering dataset derived from professional medical board exams. The corresponding knowledge graph was a subset of the Unified Medical Language System containing approximately one million entities. The second dataset is LegalBench, a collection of legal reasoning tasks including statute interpretation and case outcome prediction. The legal knowledge graph was constructed from statutes and case law citations.

**We compared our Knowledge Graph-Graph Attention Network model against three baselines:**

1. *Vanilla LLaMA-7B*: A general-purpose large language model without domain-specific fine-tuning.
2. *Fine-tuned LLaMA*: The same model fine-tuned on the training sets of the respective domains but without access to external knowledge graphs.
3. *RAG-LLaMA*: A retrieval-augmented generation approach where textual descriptions of relevant knowledge graph triples are retrieved and prepended to the input context, without using graph neural networks.

### 4.2 Implementation Details

The Graph Attention Network was implemented using the PyTorch Geometric library. We utilized a two-layer architecture with an embedding dimension of 1024 to match the hidden size of the language model. The attention mechanism employed 8 heads. The language model component was initialized with weights from LLaMA-2-7B. Training was performed using the

AdamW optimizer with a learning rate of  $2e-5$  for the language model parameters and  $1e-4$  for the Graph Attention Network parameters. A warm-up period of 1000 steps was used, followed by a linear decay of the learning rate. The training objective was the standard cross-entropy loss for next-token prediction. We trained the models on 4 NVIDIA A100 GPUs for 3 epochs. To ensure fair comparison, all baselines were trained with identical hyperparameters where applicable. The entity linking was performed using a BERT-based named entity recognition model trained on domain-specific data prior to the main experiment.

#### Code Listing 1: Graph Attention Layer Implementation

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class GraphAttentionLayer(nn.Module):
    def __init__(self, in_features, out_features, dropout, alpha, concat=True):
        super(GraphAttentionLayer, self).__init__()
        self.dropout = dropout
        self.in_features = in_features
        self.out_features = out_features
        self.alpha = alpha
        self.concat = concat

        self.W = nn.Parameter(torch.empty(size=(in_features, out_features)))
        nn.init.xavier_uniform_(self.W.data, gain=1.414)
        self.a = nn.Parameter(torch.empty(size=(2*out_features, 1)))
        nn.init.xavier_uniform_(self.a.data, gain=1.414)

        self.leakyrelu = nn.LeakyReLU(self.alpha)

    def forward(self, h, adj):
        Wh = torch.mm(h, self.W)
        a_input = self._prepare_attentional_mechanism_input(Wh)
        e = self.leakyrelu(torch.matmul(a_input, self.a).squeeze(2))

        zero_vec = -9e15*torch.ones_like(e)
        attention = torch.where(adj > 0, e, zero_vec)
        attention = F.softmax(attention, dim=1)
        attention = F.dropout(attention, self.dropout, training=self.training)
        h_prime = torch.matmul(attention, Wh)

        if self.concat:
            return F.elu(h_prime)
        else:
            return h_prime
```



```
def _prepare_attentional_mechanism_input(self, Wh):
    N = Wh.size()[0]
    Wh_repeated_in_chunks = Wh.repeat_interleave(N, dim=0)
    Wh_repeated_alternating = Wh.repeat(N, 1)
    all_combinations_matrix = torch.cat([Wh_repeated_in_chunks,
    Wh_repeated_alternating], dim=1)
    return all_combinations_matrix.view(N, N, 2 * self.out_features)
```

5. Results and Analysis

The performance of the models was evaluated using standard metrics including Accuracy, F1 Score, and Exact Match for the question-answering tasks. Additionally, we employed human evaluation for a subset of generated responses to assess coherence and factual correctness.

5.1 Quantitative Performance

The quantitative results indicate a distinct advantage for the proposed architecture. Table 1 summarizes the performance across both the MedQA and LegalBench datasets. The Vanilla LLaMA model struggled significantly with domain-specific terminology and reasoning, often defaulting to generalized but incorrect answers. Fine-tuning provided a substantial boost, particularly in learning the stylistic nuances of the domains, but still suffered from hallucinations on complex queries involving rare entities.

Table 1: Experimental Results Comparison

Model	MedQA Accuracy (%)	MedQA F1 Score	LegalBench Accuracy (%)	LegalBench F1 Score
Vanilla LLaMA-7B	34.2	31.5	41.8	38.2
Fine-tuned LLaMA	48.6	46.1	55.3	52.7
RAG-LLaMA	54.1	51.8	61.2	59.4
KG-GAT (Ours)	62.8	60.4	68.7	66.1

The Retrieval-Augmented Generation baseline showed improvement over simple fine-tuning, validating the hypothesis that external knowledge is crucial. However, our Knowledge Graph-Graph Attention Network (KG-GAT) model outperformed the RAG baseline by a significant margin—approximately 8.7% in accuracy on MedQA and 7.5% on LegalBench. This suggests that the structured integration of knowledge via Graph Attention Networks is more effective than unstructured text retrieval. The attention mechanism allows the model to synthesize information from multiple connected nodes, effectively performing multi-hop reasoning that is difficult to achieve with simple text retrieval [14]. Furthermore, the improvement in F1 scores indicates that our model captures the precise terminology required in these domains more accurately [15].

5.2 Ablation Studies

To understand the contribution of individual components, we conducted ablation studies by removing specific parts of the architecture. We tested a variant where the Graph Attention Network was replaced with a static Graph Convolutional Network (GCN), and a variant where the graph injection was performed only at the input layer rather than deep within the network.

Table 2: Ablation Study Results

Configuration	MedQA Accuracy (%)	Accuracy Delta	LegalBench Accuracy (%)	Delta
KG-GAT (Full Model)	62.8	-	68.7	-
w/o Attention (GCN)	58.3	-4.5	64.1	-4.6
w/o Deep Injection	56.9	-5.9	63.5	-5.2
w/o Graph Structure	54.1	-8.7	61.2	-7.5

The results in Table 2 demonstrate the critical role of the attention mechanism. Replacing the Graph Attention Network with a Graph Convolutional Network resulted in a performance drop, confirming that treating all neighbors equally introduces noise that hampers model performance. The drop was even more pronounced when deep injection was removed, highlighting the importance of fusing knowledge at higher levels of abstraction. The "w/o Graph Structure" variant essentially replicates the RAG baseline, confirming that the topological information contained in the graph edges provides valuable signal beyond the content of the nodes themselves [16].

6. Discussion

The experimental findings underscore the transformative potential of integrating structured knowledge graphs into deep learning pipelines via attention mechanisms. The superior performance of the KG-GAT model suggests that the "black box" nature of large language models can be effectively illuminated by the structured logic of symbolic AI.

6.1 Implications for Domain Specialization

One of the most significant implications of this research is the reduction in data requirements for domain adaptation. Traditional fine-tuning requires massive corpora of domain-specific text to implicitly teach the model facts. Our approach decouples the reasoning engine (the Language Model) from the knowledge base (the Knowledge Graph). This means that a model can be updated with new facts simply by updating the Knowledge Graph, without the need for expensive retraining. This is particularly valuable in fields like medicine or law, where knowledge is constantly evolving.



Figure 2: Efficiency Chart



Figure 2 illustrates this efficiency. Even with a fraction of the training data, the KG-GAT model achieves performance comparable to a fully fine-tuned model, provided the underlying Knowledge Graph is comprehensive. This suggests a paradigm shift where the focus of domain adaptation moves from collecting unstructured text to curating high-quality structured data.

## 6.2 Limitations and Challenges

Despite the promising results, several limitations persist. First, the performance of the system is heavily dependent on the quality and completeness of the underlying Knowledge Graph. If the graph contains errors or lacks coverage, the model's performance will suffer, a phenomenon known as error propagation. Second, the computational overhead of processing graphs with Graph Attention Networks is non-negligible. While more efficient than processing equivalent amounts of text, the graph retrieval and attention mechanisms add latency to the inference process [17]. Furthermore, the entity linking step acts as a bottleneck. If entities in the input text are not correctly identified or linked to the graph, the subsequent injection mechanism fails to provide relevant context. Future work must address robustness against noisy entity linking and explore end-to-end differentiable linking strategies.

## 6.3 Future Directions

Future research should explore the integration of dynamic Knowledge Graphs that can evolve during the conversation. Additionally, extending this architecture to multi-modal graphs that include images or numerical data could further enhance its applicability in domains like radiology or finance. Investigating the interpretability of the attention weights in the Graph Attention Network could also provide valuable insights into the reasoning process of the model, offering a form of "explainable AI" that is highly sought after in regulated industries [18].

## Conclusion

In this paper, we presented a comprehensive framework for enhancing domain-specific Large Language Models using Knowledge Graph Injection and Graph Attention Networks. By explicitly modeling the relationships between entities and dynamically attending to relevant structural context, our approach significantly reduces hallucination and improves factual accuracy in specialized domains. The empirical results on biomedical and legal datasets confirm that the synergy between symbolic knowledge representation and neural probabilistic modeling is a powerful direction for the future of artificial intelligence. As the field moves toward more robust and reliable systems, architectures that respect and utilize structured knowledge will likely become the standard for high-stakes applications. The proposed KG-GAT framework represents a significant step toward bridging the gap between the statistical fluency of language models and the rigorous precision required by expert domains.

## References

- [1] Zhang, T. (2025). A Neuro-Symbolic and Blockchain-Enhanced Multi-Agent Framework for Fair and Consistent Cross-Regulatory Audit Intelligence.
- [2] Zhang, W., Zhang, C., Luo, Z., Ma, J., Yuan, W., Gu, C., & Feng, C. (2025). SemanticForge: Repository-Level Code Generation through Semantic Knowledge Graphs and Constraint Satisfaction. arXiv preprint arXiv:2511.07584.
- [3] Liu, J., Kong, Z., Zhao, P., Yang, C., Shen, X., Tang, H., ... & Wang, Y. (2025, April). Toward adaptive large language models structured pruning via hybrid-grained weight importance assessment.

- In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 39, No. 18, pp. 18879-18887).
- [4] Li, J., & Cappelleri, D. J. (2024). Sim-grasp: Learning 6-dof grasp policies for cluttered environments using a synthetic benchmark. *IEEE Robotics and Automation Letters*.
  - [5] Yi, X. (2025, October). Compliance-by-Design Micro-Licensing for AI-Generated Content in Social Commerce Using C2PA Content Credentials and W3C ODRL Policies. In 2025 7th International Conference on Machine Learning, Big Data and Business Intelligence (MLBDBI) (pp. 204-208). IEEE.
  - [6] Liu, F., Tian, J., Miranda-Moreno, L., & Sun, L. (2023). Adversarial danger identification on temporally dynamic graphs. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4), 4744-4755.
  - [7] Zhao, J. Analysis of working women's perceptions of state-regulated family planning policy: China as a case study (Doctoral dissertation, Loughborough University).
  - [8] Zhang, W., Zhang, C., Gu, C., Kou, J., Yuan, H., Fang, X., ... & Fang, Y. (2024, October). Hallucination in Large Language Models: From Mechanistic Understanding to Novel Control Frameworks. In 2024 7th International Conference on Universal Village (UV) (pp. 1-36). IEEE.
  - [9] Sun, Q., Zhao, X., & Lin, X. (2025, September). Design of a Hardware-Software Co-designed Real-Time Machine Learning System for Big Data Streams. In Proceedings of the 2nd International Symposium on Integrated Circuit Design and Integrated Systems (pp. 265-271).
  - [10] Yang, Y., Lin, Z., & Wei, L. (2025). ACE-Sync: An Adaptive Cloud-Edge Synchronization Framework for Communication-Efficient Large-Scale Distributed Model Training. *arXiv preprint arXiv:2512.18127*.
  - [11] Li, J., & Cappelleri, D. J. (2023). Sim-suction: Learning a suction grasp policy for cluttered environments using a synthetic benchmark. *IEEE Transactions on Robotics*, 40, 316-331.
  - [12] Zhou, Z., Zhao, C., Li, X., Zhang, H., & Chang, R. (2025, July). Diverse Stacking Ensemble for Attributing LLM Outputs via Relational Reasoning. In 2025 8th International Conference on Computer Information Science and Application Technology (CISAT) (pp. 1089-1092). IEEE.
  - [13] Liu, F., Jiang, S., Miranda-Moreno, L., Choi, S., & Sun, L. (2024). Adversarial vulnerabilities in large language models for time series forecasting. *arXiv preprint arXiv:2412.08099*.
  - [14] Hu, Z., Chen, X., & Hu, J. (2025). Emotion-Driven Personalized Recommendation for AI-Generated Content Using Multi-Modal Sentiment and Intent Analysis. *arXiv preprint arXiv:2512.10963*.
  - [15] Liu, S., Du, H., & Wang, S. (2025). Adaptive Cache Pollution Control for Large Language Model Inference Workloads Using Temporal CNN-Based Prediction and Priority-Aware Replacement. *arXiv preprint arXiv:2512.14151*.
  - [16] Bai, Z., & Chen, K. (2025, September). Study on Adaptive Optimisation Method for AI Generated Code Performance Based on Reinforcement Learning. In Proceedings of the 2nd International Symposium on Integrated Circuit Design and Integrated Systems (pp. 185-190).
  - [17] Zhang, H., Zhao, S., Zhou, Z., Zhang, W., & Meng, Y. (2025, September). Domain-Specific RAG with Semantic Normalization and Contrastive Feedback for Document Question Answering. In 2025 7th International Conference on Internet of Things, Automation and Artificial Intelligence (IoTAAI) (pp. 750-753). IEEE.
  - [18] Liu, F., & Liu, C. (2018, June). Towards accurate and high-speed spiking neuromorphic systems with data quantization-aware deep networks. In Proceedings of the 55th Annual Design Automation Conference (pp. 1-6).