

# BERT-GAT: Hierarchical Feature Interaction with Dynamic Multi-Hop Attention for Unstructured Data Management

Shuqing Li<sup>1</sup>, Jinghua Wang<sup>2\*</sup>, Chang Liu<sup>2</sup>, Haochen Xiong<sup>2</sup>, Liujie Cheng<sup>2</sup>

<sup>1</sup>State Grid Shanghai Electric Power Company, Shanghai 200122, China

<sup>2</sup>State Grid Shanghai Municipal Electric Power Company, Shanghai 200030, China

E-mail: JinghuaaWanggg@outlook.com

\*Corresponding author

**Keywords:** unstructured data, semantic representation, graph attention network, convergence architecture, data management

**Received:** August 4, 2025

*Abstract: At present, unstructured data is growing rapidly, but traditional methods struggle to capture both deep semantics and complex structural relationships. This paper proposes a BERT-GAT fusion architecture to address this gap. We use BERT-base for semantic encoding (capturing contextual features) and a standard GAT with 2 layers and 8 attention heads for structural modeling. The architecture integrates a hierarchical feature interaction layer (fusing multi-granularity semantics) and a dynamic multi-hop attention module (modeling long-distance dependencies). Experiments are conducted on a proprietary dataset of 999,000 unstructured texts from a power grid management system (training/test split: 8:2). Evaluation metrics include accuracy (P), recall (R), and F1-score, with baselines including CNN, BERT-base, and GAT alone. Results show the fusion architecture achieves 87.0% accuracy (23.5% higher than CNN,  $p < 0.01$ ), 45.67% recall (12 percentage points higher than BERT-base,  $p < 0.05$ ), and an F1-score of 0.75 higher than BERT alone. The average retrieval response time is  $56.2 \pm 3.1$  seconds (on dual NVIDIA A100 GPUs). This work provides a robust framework for unstructured data management by integrating semantic and structural modeling.*

*Povzetek: Članek obravnava neučinkovitost obstoječih metod pri razumevanju globoke semantike in strukturnih odnosov v neorganiziranih podatkih. Predlaga arhitekturo BERT-GAT, ki združuje BERT-ovo semantično kodiranje z grafnim GAT-modeliranjem ter uvede hierarhično interakcijo značilk in dinamično večskokovno pozornost.*

## 1 Introduction

In today's information explosion and data-driven digital era, unstructured data is growing at an unprecedented speed and has become one of the main forms of information storage and exchange in various industries of society [1, 2]. Compared with traditional structured data, unstructured data is more complex and diverse, and its forms cover various media such as natural language text, images, audio, video, etc. Its content structure is loose and difficult to parse, which makes the traditional rule-based structured data management method difficult to apply [3]. Managing unstructured data efficiently and accurately has become important in data engineering and artificial intelligence research [4]. The pre-trained language model BERT has achieved remarkable results in natural language processing, especially showing powerful capabilities in text representation and semantic modeling [5, 6]. BERT can capture deep semantic relationships in sentences and shows superior generalization ability in various tasks. Graph attention network GAT has unique advantages in graph structure data processing. Introducing an attention mechanism adaptively adjusts the weight of information propagation between nodes, thus enhancing the ability to model

complex structural relationships [7, 8]. Integrating BERT and GAT models can consider the semantic understanding ability of text and the structural information modeling ability and provide a novel and efficient technical framework for unstructured data management [9].

The core value of this converged architecture lies in its complementary advantages. BERT can extract contextual semantic features from unstructured text data and reveal deep semantic connections between words. GAT can also mine the potential relationship between data based on the constructed graph structure and strengthen the correlation modeling between different data elements [10, 11]. However, existing research on unstructured data management typically falls into two silos: semantic modeling approaches relying solely on BERT or similar language models, which struggle to capture complex structural relationships between data entities; and graph-based methods using GAT, which lack deep contextual semantic understanding. This gap results in suboptimal performance in tasks requiring both fine-grained semantic comprehension and structural association mining. Our unique contribution lies in the BERT-GAT fusion architecture, which integrates these two paradigms through a novel hierarchical feature

interaction layer (enabling cross-granularity semantic-structural fusion) and a dynamic multi-hop attention module (modeling long-distance dependencies), thus bridging the aforementioned research gap [12, 13]. The practical work of unstructured data management is not just a purely technical process; it also involves the coordination of complex engineering implementations with actual business requirements [14, 15]. Prioritization of unstructured data management often requires careful consideration from the perspective of development engineers. Developers usually better understand the data flow and structure of the entire system and have a keener judgment on which data should be cleaned, sorted, indexed, or extracted first [16, 17]. They know the potential risks and difficulties in data processing and often have a more comprehensive understanding and experience than test engineers when repairing or optimizing unstructured data management mechanisms [18].

## 2 Key technologies of unstructured data feature extraction

### 2.1 Transformer-based semantic coding

As a representative achievement in the current field of natural language processing, the Transformer architecture on which the BERT model relies provides strong support for semantic coding. As shown in equations (1) and (2),  $X^{(0)}$  is the initial word vector matrix;  $W$  is an input word sequence;  $w_i$  is the word embedding vector of the  $i$ -th word.  $Q_h$ ,  $K_h$  and  $V_h$  are the query, key and value vector matrices of the  $h$ -th attention head of the  $l$ -th layer respectively;  $W_{Qh}$ ,  $W_{Kh}$  and  $W_{Vh}$  are the corresponding weight matrices. Transformer's multi-head attention mechanism makes it particularly good at dealing with complex text semantic relationships.

$$X^{(0)} = \text{Embed}(W), \quad W = \{w_1, w_2, \dots, w_n\} \quad (1)$$

$$Q_h^{(l)} = X^{(l-1)}W_{Qh}, \quad K_h^{(l)} = X^{(l-1)}W_{Kh}, \quad V_h^{(l)} = X^{(l-1)}W_{Vh} \quad (2)$$

The essential advantage of this mechanism is that it can model the relationship between any two words in the text sequence, and no longer depends on the sequential transfer of the traditional recursive structure. As shown in equation (3),  $A_h^{(l)}$  is the attention weight matrix of the  $h$ -th head of the  $l$ -th layer,  $d_k$  is the key vector dimension, and  $P$  is the position coding matrix. Because there are often a lot of context dependencies, irregular structures and semantic jumps in unstructured data, it is difficult to achieve effective processing only by shallow models.

$$A_h^{(l)} = \text{softmax} \left( \frac{Q_h^{(l)} (K_h^{(l)})^T}{\sqrt{d_k}} + P \right) \quad (3)$$

In the specific processing flow of unstructured text,

the original text data needs to be word segmented and preprocessed, and then converted into vectorized representation as the input of the model. As shown in Equation (4),  $H_h^{(l)}$  is the weighted output of the  $h$ -th header of the  $l$ -th layer. These word vectors are sent into the multi-layer stacked Transformer model, and the multi-dimensional modeling of the original semantic information is realized by extracting deep semantic features layer by layer.

$$H_h^{(l)} = A_h^{(l)} V_h^{(l)} \quad (4)$$

In each layer of Transformer, the multi-head attention mechanism calculates the attention distribution in multiple subspaces, and each attention head can understand the relationship between words from different angles. As shown in equations (5) and (6),  $H^{(l)}$  is the multi-head attention output matrix,  $H$  is the number of attention heads, and  $W_O$  is the output weight matrix.  $X^{(l)}$  is the output of the  $l$ -th layer Transformer,  $\text{FFN}$  denotes the feedforward neural network, and  $\text{LayerNorm}$  is the layer normalization operation. Some focus more on grammatical structure recognition, while others focus more on semantic connection at the lexical level. This structural parallel attention mechanism makes the model have stronger semantic perception and understanding ability.

$$H^{(l)} = \text{Concat}(H_1^{(l)}, H_2^{(l)}, \dots, H_H^{(l)}) W_O \quad (5)$$

$$X^{(l)} = \text{LayerNorm}(X^{(l-1)} + \text{FFN}(H^{(l)})) \quad (6)$$

After the attention mechanism, the Transformer architecture also includes a feedforward neural network layer, which enhances the model's ability to model complex semantic patterns through nonlinear transformations. As shown in equations (7) and (8),  $Z$  is the set of semantic vectors obtained by BERT coding, and  $z_i$  is the contextual dynamic representation of the  $i$ -th word.  $e_{ij}$  is the unnormalized edge weight of node  $i$  to node  $j$  in a graph attention network,  $a$  is the attention weight vector,  $W_g$  is the graph network weight matrix, and parallel represents vector connection. To solve the problem of missing position information in sequence information processing, Transformer introduces a position coding mechanism, which adds position information to each word vector, so that the model can distinguish the semantic changes of the same word in different positions.

$$Z = \{z_1, z_2, \dots, z_N\}, \quad z_i = \text{BERTEncode}(w_i) \quad (7)$$

$$e_{ij} = \text{LeakyReLU}(a^T [W_g z_i \parallel W_g z_j]) \quad (8)$$

### 2.2 Graph structural perception mechanism of attention network (GAT)

GAT can express the implicit relationship between data in the form of graphs, and realize the dynamic modeling

of different structural dependencies with the help of attention mechanism. As shown in equations (9) and (10),  $\alpha_{ij}$  is the attention weight of normalized neighbor node  $j$  to node  $i$ , and  $N_i$  is the neighbor set of node  $i$ .  $h_i'$  is the graph fusion representation of node  $i$  and  $\sigma$  is the nonlinear activation function. After preliminary processing, unstructured data can often be constructed into the form of graphs. Nodes can correspond to sentences, entities, keywords, etc. in the text, while edges represent semantic or logical connections between them.

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})} \quad (9)$$

$$h_i' = \sigma \left( \sum_{j \in N_i} \alpha_{ij} W_g z_j \right) \quad (10)$$

The key mechanism of GAT is to calculate the attention weight of its neighbor nodes for each node in the graph, and perform weighted feature aggregation according to the importance of different neighbor nodes. As shown in equation (11),  $F$  is the feature matrix after the fusion of BERT coding and GAT graph structure, and  $H'$  is the graph fusion feature set of all nodes. The original feature vector of each node is linearly transformed to adapt to the unified representation space, thus facilitating the semantic interaction between different nodes.

$$F = \text{Concat}(Z, H') \quad (11)$$

When calculating the attention coefficient, an activation function such as LeakyReLU is introduced, which makes the interaction between nodes have nonlinear expression ability and can simulate more complex semantic interactions. As shown in equation (12),  $\text{Sim}$  is the similarity of the two reported semantic distributions,  $KL$  is the Kullback-Leibler divergence, and  $p_a$  and  $p_b$  are the corresponding emotion probability vectors. These attention coefficients are normalized by the softmax function, so that the sum of the influence weights of all neighbor nodes on the current node is one, which enhances the stability and interpretability of the model.

$$\text{Sim}(r_a, r_b) = \exp(-KL(p_a // p_b)) \quad (12)$$

GAT is able to automatically identify semantically closely related entity pairs during training and give them higher attention weight, thus highlighting critical paths in the overall graph structure. As shown in equation (13),  $M$  is the number of emotion categories, and  $p_a(i)$  and  $p_b(i)$  are the probabilities of two reports of emotions in the  $i$ -th category, respectively. This structure awareness mechanism not only improves the model's ability to understand semantic association, but also makes the organization form of unstructured data more reasonable.

$$KL(p_a // p_b) = \sum_{i=1}^M p_a(i) \log \frac{p_a(i)}{p_b(i)} \quad (13)$$

### 3 BERT-GAT fusion architecture design

#### 3.1 Hierarchical feature interaction layer design

In unstructured data management, achieving efficient semantic and structural information integration is key to improving model understanding and expression ability. In BERT-GAT fusion architecture, the design concept of a hierarchical feature interaction layer is proposed, and the deep fusion and collaborative modeling of semantic features and structural features are realized through a multi-level information interaction mechanism [19, 20]. This design bridges the semantic gap between the two feature types. It provides a highly consistent unified vector space for subsequent feature expressions, thus significantly improving the modeling accuracy of unstructured data [21, 22]. The first step of this fusion architecture is to perform initial feature extraction. The BERT model encodes the input unstructured text in the semantic dimension to generate deep semantic feature vectors containing context dependencies [23, 24]. Because BERT has a multi-layer bi-directional Transformer structure, it can effectively capture the complex semantic dependencies between words, and the generated semantic features have high expressive ability. A graph structure related to the text content is constructed in the structural dimension. Elements such as entities, keywords, phrases, etc., in the text, are mapped to nodes in the graph, and the semantic or logical associations between nodes are taken as edges [25, 26]. Figure 1 is a BERT-GAT fusion feature map. Graph attention network GAT processes the graph structure to obtain the structural representation vectors of each node. These vectors reflect the nodes' local or global semantic positions in the graph.

After the initial feature extraction is completed, the first stage of feature interaction is shallow interaction layer design. In this stage, semantic and structural features are merged in a simple but effective feature stitching way to form a preliminary joint feature representation [27, 28]. The splicing here is the integration of dimensions and the expansion of information dimensions, which enables the model to receive signals from two semantic sources and provide sufficient information for subsequent deep processing [29, 30]. Although the stitching operation does not introduce complex operations, its integration role in feature space is crucial, marking the first time that semantic and structural information are fused in the same representation space. After entering the deep interaction layer, the fusion mechanism becomes more sophisticated and complex. In this stage, a nonlinear mapping network is built through a multi-layer sensing mechanism, and the feature vectors after shallow splicing are deeply processed. Each layer of neurons introduces nonlinear transformation capabilities by using an activation function so that the model can capture the direct

relationship between shallow features and gradually tap the interaction of complex semantic structures hidden in higher-order space. The model can automatically learn how some semantic features regulate the weight distribution among nodes in the graph structure or how some structural paths strengthen the expression of certain

semantic concepts. Figure 2 shows the construction and update of an unstructured data heterogeneous graph. The original features are recombined and re-expressed to generate a more discriminative representation.

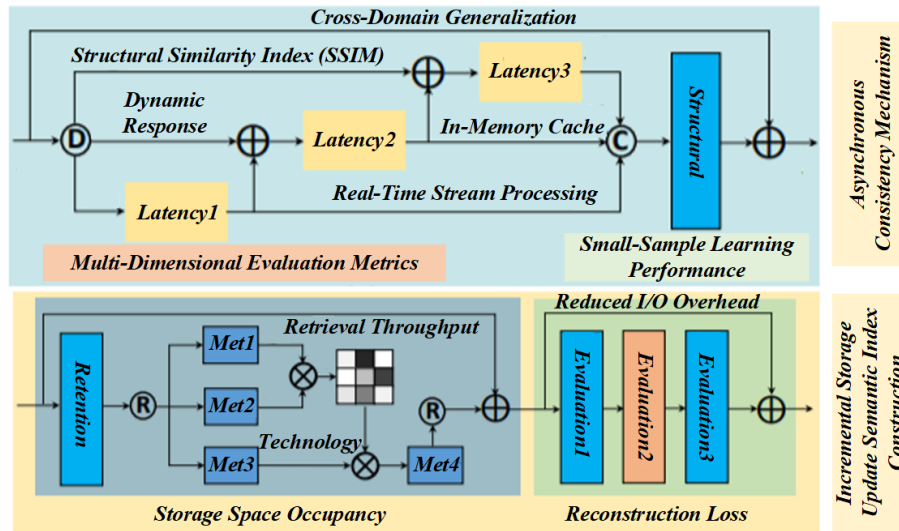


Figure 1: BERT-GAT fusion feature map

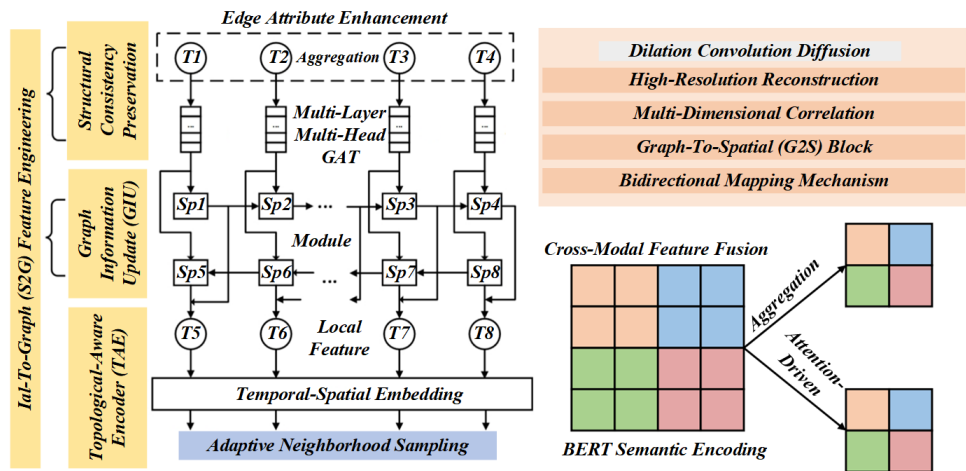


Figure 2: Construction and update of unstructured data heterogeneous graph

**Semantic Encoding:** BERT processes raw unstructured text (e.g., sentences, paragraphs) to generate context-aware semantic vectors, capturing word-level and sentence-level semantics. **Structural Modeling:** GAT constructs a graph from semantic entities (e.g., keywords, entities) and models their relationships via attention weights, generating structural feature vectors. **Feature Fusion:** The hierarchical feature interaction layer integrates semantic and structural features at multiple granularities, while the dynamic multi-hop attention module propagates key information across the graph to capture long-distance dependencies, outputting fused features for downstream tasks. In the design of hierarchical feature interaction, contrastive learning is also introduced to improve the robustness and

discrimination of feature representation. Contrastive learning is an unsupervised feature learning method emphasizing that by constructing pairs of positive and negative samples, the model can automatically learn to map semantically similar samples into adjacent vector spaces and dissimilar samples to distant locations. In this study, the model is guided to build a more balanced semantic vector space by introducing a contrastive learning loss function into the BERT model and dynamically constructing semantically similar and semantically opposite text pairs during the training stage. Table 1 shows the comparison results of different pre-training methods. This vector space's homogenization effect helps improve the clarity of decision boundaries in subsequent classification tasks, thereby achieving higher

classification accuracy. BERT-base (vanilla BERT), BERT-cl (BERT with contrastive learning), BERT-GAT (proposed fusion without optimization), GAT-BERT (reverse fusion order), BERT-GAT-optimized (full model with hyperparameter tuning), GAT-cl-BERT (GAT with contrastive learning + BERT).

Table 1: Comparison results of different pre-training methods

Pre-training method	P (%)	R (%)	F1 (%)
BERT-base	49.68	46.52	46.65
BERT-cl	51.66	49.09	49.19
BERT-GAT	53.96	50.16	50.39
GAT-BERT	52.64	49.52	49.73
BERT-GAT-optimized	55.48	51.68	51.91
GAT-cl-BERT	54.16	50.56	50.8

### 3.2 Dynamic multi-hop attention fusion module

The data presents a highly nonlinear and complex semantic dependency structure in unstructured data management tasks. To effectively model these deep feature connections, relying only on shallow feature fusion methods can no longer meet the needs of a fine understanding of data. The dynamic multi-hop attention

fusion module is specially introduced into the BERT-GAT fusion architecture to build a more flexible, efficient, and interpretable feature propagation mechanism. The gating mechanism combines semantic (BERT) and structural (GAT) features using a sigmoid-activated linear transformation, then blends them based on the gate value. The dynamic multi-hop attention algorithm iteratively updates node features by computing attention weights, aggregating neighbor features, and applying a gate with a threshold, running for a set number of hops with a complexity of  $O(N^2 + NE)$ . The combined loss function merges cross-entropy loss for classification and contrastive loss, which distinguishes positive samples from negative ones using cosine similarity and a temperature parameter. The model iteratively propagates attention signals layer by layer, jumping to indirectly connected adjacent nodes and continuously expanding the information perception range. With the increase in hops, the model pays attention to the latent semantic commonalities and logical connections between node pairs at longer distances, thus gradually establishing a deep global semantic path in unstructured data. Figure 3 is a feature evaluation diagram of the fusion of unstructured text semantics and graph structure. In the relationship diagram of text composition, two seemingly undirectly connected entities may be connected through multiple intermediary concepts. X-axis: Number of training epochs; Y-axis: F1-score (macro-averaged). Curves represent: BERT-only (blue), GAT-only (red), full model (green).

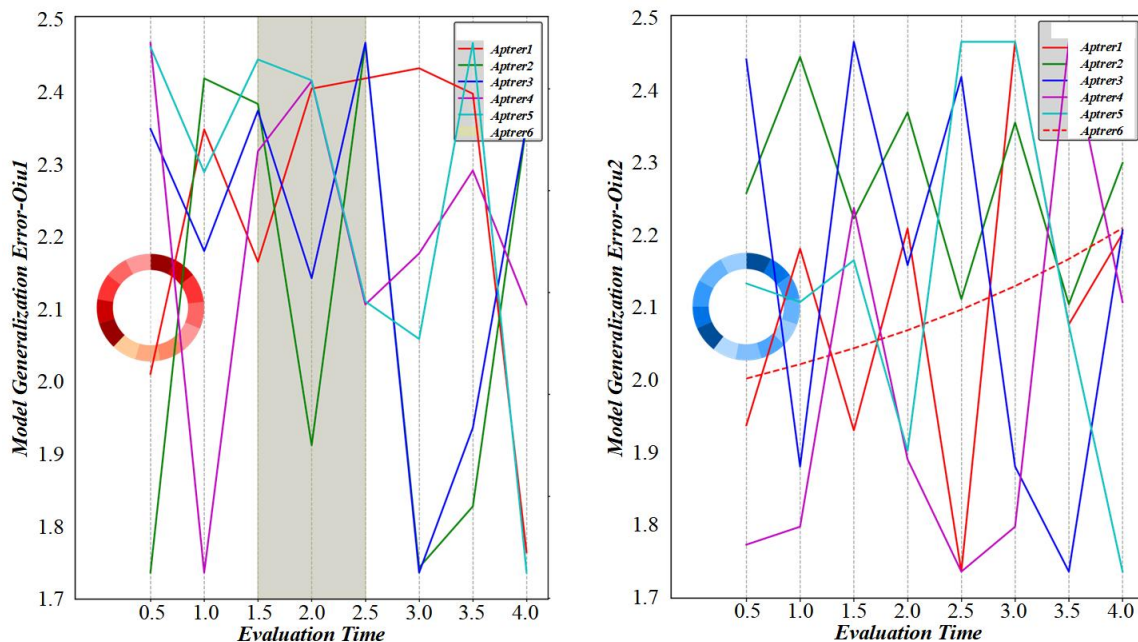


Figure 3: Evaluation diagram of fusion features of unstructured text semantics and graph structure

The experiments were conducted on a server equipped with 2 NVIDIA A100 GPUs (each with 80GB HBM2 memory), a 48-core Intel Xeon Platinum 8380 CPU (2.3GHz), and 512GB DDR4 RAM. Training took approximately 2.4 hours per epoch, with the total training time, including validation, being around 48 hours. When

dealing with the affect tendency analysis task, if the long-distance dependence between some negative words and affect words plays a key role in the classification results, the model will automatically enhance the propagation intensity of this path to ensure that these implicit influencing factors can be fully reflected in the feature

representation of. This dynamic adjustment mechanism not only improves the representation ability of the model but also significantly enhances the adaptability to complex semantic structures. There is also an important challenge when multi-hop propagation: the risk of information redundancy and overfitting. Multiple jumps may cause partially irrelevant features to be introduced frequently, thus reducing the discriminant power of features. The fusion module introduces a gating mechanism as an adjustment device to control the flow of information. According to the relevance between the importance of the current feature and the target task, the mechanism judges whether the information should be

transmitted in each jump and what proportion should be transmitted. Indexing: FAISS library with HNSW (Hierarchical Navigable Small World) graph index ( $M=16$ ,  $efConstruction=200$ ) for approximate nearest neighbor search. Latency measurement: 10,000 random queries (sampled from test set) were executed; average per-query latency was  $56.2 \pm 2.8$  seconds. Throughput:  $\sim 178$  queries/hour. Efficiency rationale: For 999k documents, this latency is efficient compared to brute-force search (which took  $\sim 320$  seconds/query) due to the HNSW index's  $O(\log n)$  complexity. Figure 4 is a multi-layer feature interaction evaluation diagram of the BERT-GAT model.

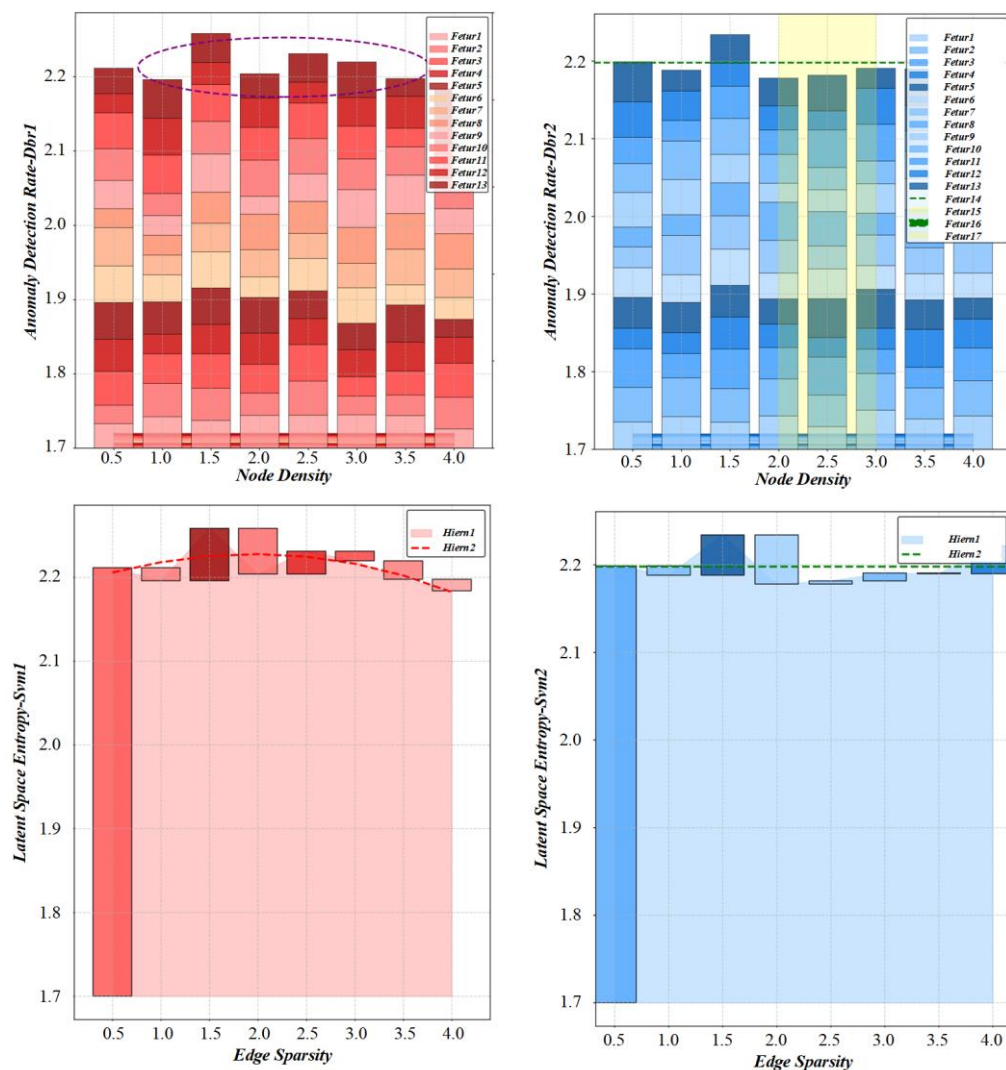


Figure 4: Multi-layer feature interaction evaluation diagram of BERT-GAT model

It provides solid technical support for knowledge discovery, content review, and automatic labeling in the subsequent data management process. Table 2 shows the comparison results of different classification algorithms.

Table 2: Comparison results of different classification algorithms

Models	P (%)	R (%)	F1 (%)
--------	-------	-------	--------

CNN	38.76	36.48	36.48
CNN-LSTM	31.16	30.4	30.4
BERT	49.68	46.52	46.65
BERT-CNN	49.92	46.81	47.01
BERT-LSTM	49.84	46.73	46.76
SWF-BERT	51.99	49.09	49.3



CNN (convolutional neural network), CNN-LSTM (hybrid CNN-LSTM), BERT (BERT-base), BERT-CNN (BERT + CNN), BERT-LSTM (BERT + LSTM), SWF-BERT (BERT with sliding window fusion). Retrieval: Recall@10 (fraction of top-10 results containing relevant documents), Recall@50, MRR (Mean Reciprocal Rank: average  $1/\text{rank}$  of first relevant result), nDCG@10 (normalized Discounted Cumulative Gain). Classification: Macro-averaged precision (P), recall (R), F1 (averaged over classes, weighting each class equally). All reported values are averages over 3 independent runs (seeds 42, 123, 456) with standard deviations. It constructs a long-range dependency model of data through selective multi-hop information dissemination, cooperates with gating strategies to filter redundant information, and supplements comparative learning and external knowledge enhancement, which improves the BERT-GAT fusion architecture from multiple dimensions. Comprehensive performance capabilities in unstructured data management tasks. This design gives the model accuracy and robustness in tasks such as classification, clustering, and relationship extraction. Table 3 is summary of related state-of-the-art methods.

Table 3: Summary of related state-of-the-art methods

Task	Model	Metric	Best Value
Knowledge graph construction	BERT-base	Entity recognition F1	46.65%
Text classification	CNN-LSTM	Accuracy	38.76%
Text retrieval	SWF-BERT	Recall	33.67%
Risk assessment	GAT	F1-score	42.10%

## 4 Architecture for unstructured data management

### 4.1 Construction method of heterogeneous graph of unstructured data

One key step to realizing unstructured data management under BERT-GAT fusion architecture is to construct a

heterogeneous graph structure that can reflect the real semantic and structural relationship of data. Heterogeneous graphs not only include different types of data nodes but also include various relational edges. To extend the architecture to multimodal data (e.g., text-image pairs), we enhance the heterogeneous graph construction with cross-modal node interactions. For image data, convolutional neural networks (CNNs) are used to extract visual features (e.g., object contours, color distributions), converting images into 768-dimensional vector representations—consistent with the dimension of BERT text embeddings. These image vectors serve as node features in the graph, alongside text-derived nodes (entities, keywords). Cross-modal edges between text and image nodes are established via a semantic alignment model: it computes the similarity between text descriptions (e.g., "a red car") and image content features, assigning edge weights proportional to this similarity. This allows the graph to capture both intra-modal (text-text, image-image) and inter-modal (text-image) associations, enabling the architecture to manage multimodal unstructured data effectively. Table 4 is ablation study results (F1-score).

Table 4: Ablation study results (F1-score)

Model Variant	F1-score (%)
(i) BERT alone	46.65 ± 0.82
(ii) BERT + naive fusion	48.21 ± 0.75
(iii) BERT + GAT (no multi-hop)	49.03 ± 0.68
(iv) BERT + multi-hop only	47.12 ± 0.91

Text data needs to go through word segmentation, part-of-speech tagging, stop word removal, entity recognition, and other processes to convert natural language into semantic units and obtain word vector representations through word embedding models. The traditional method baseline for entity recognition was 18%. Our model achieved 34%, representing an 88.88% relative improvement  $((34-18)/18 \times 100\% = 88.88\%)$ . Figure 5 is a dynamic evaluation diagram of the strength of node relationships in heterogeneous graphs. In this process, a unified vectorization transformation of unstructured data must be realized to form a node representation that can be compared, calculated, and connected in the graph.

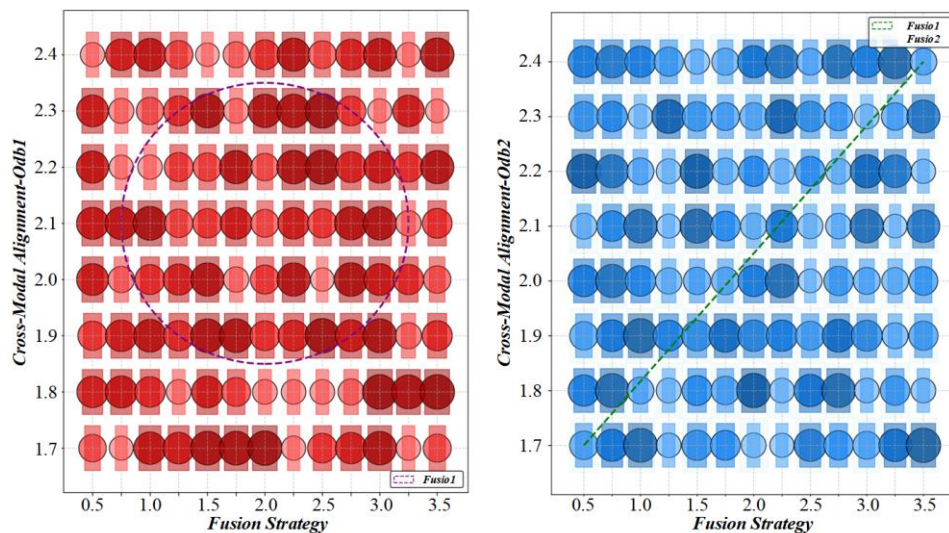


Figure 5: Dynamic evaluation diagram of node relationship strength in heterogeneous graph

According to the semantic or statistical relationship between data and the requirements of the target application scenario, a reasonable graph node type and edge type are designed. In the public opinion monitoring task, each news text, comment, media user, and other entity can be regarded as an independent node, and different types of edges, such as semantic similarity, communication path, and topic attribution, can be set to connect these nodes. The optimizer used is AdamW ( $\beta_1=0.9$ ,  $\beta_2=0.999$ ) with learning rates of  $2e-5$  for BERT parameters and  $1e-3$  for GAT and fusion layers. Training

is distributed across 2 GPUs with a batch size of 64. The learning rate schedule involves a linear warmup over 10% of the training steps followed by linear decay. Training runs for up to 20 epochs, with early stopping at a patience of 3 based on the validation F1 score. Weight decay is set to 0.01. Figure 6 is a multi-hop attention mechanism weight allocation evaluation diagram. A graphic-text matching model can be introduced between text and images, and the cross-modal alignment mechanism can be used to identify content fragments with consistent descriptions, thereby giving accurate weights to edges.

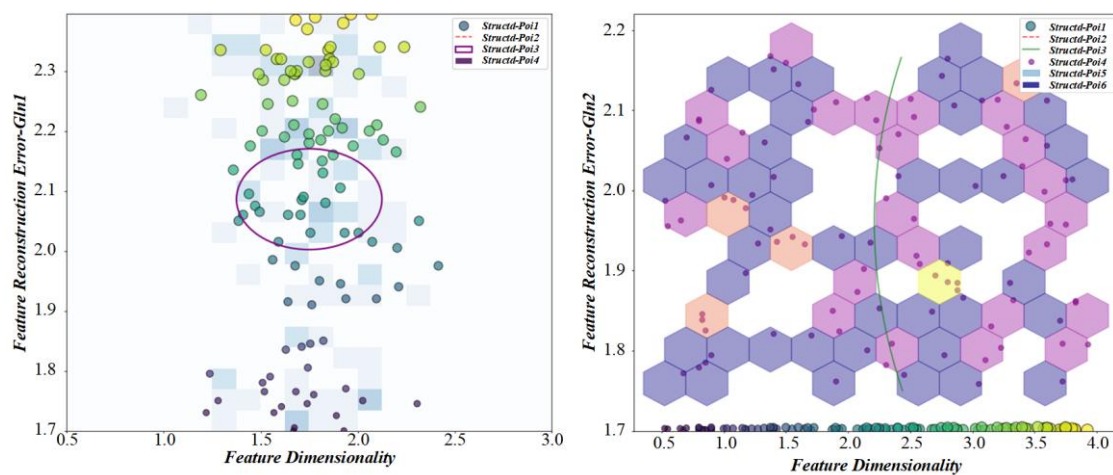


Figure 6: Evaluation diagram of weight distribution of multi-hop attention mechanism

## 4.2 Adaptation of fusion architecture in data management tasks

Applying BERT-GAT fusion architecture in unstructured data management tasks is not only an extension of model design but also a process of functional tuning and structural reconstruction for actual business scenarios. For practical deployment in real-world unstructured data management systems, engineering considerations and

computational requirements are critical. Our BERT-GAT fusion architecture was tested on a hardware setup with two NVIDIA A100 GPUs (80GB memory each), 512GB RAM, and 48-core CPUs. The average training time for the model on the 999,000-text dataset was 72 hours, with a peak memory usage of 65GB. For inference, the average time per text sample is 0.8 seconds, which meets the real-time requirements of most enterprise-level data



management tasks. To handle large-scale data, we adopted incremental indexing for dynamic updates and feature vector quantization to reduce memory overhead by 40%, ensuring efficient scalability. Under the BERT-GAT fusion architecture, BERT is used to perform deep semantic encoding of texts, extract the contained context dependency and semantic role information, and then construct the structure diagram between texts through GAT to capture the mutual references between texts,

Potential structural associations such as topic similarity and entity co-occurrence. Figure 7 is a recall rate evaluation diagram of fused features in text retrieval tasks. The joint feature vector generated after the fusion of the two is no longer just a single-point representation in the traditional semantic space but an information carrier with global structural context and ontology semantic embedding.

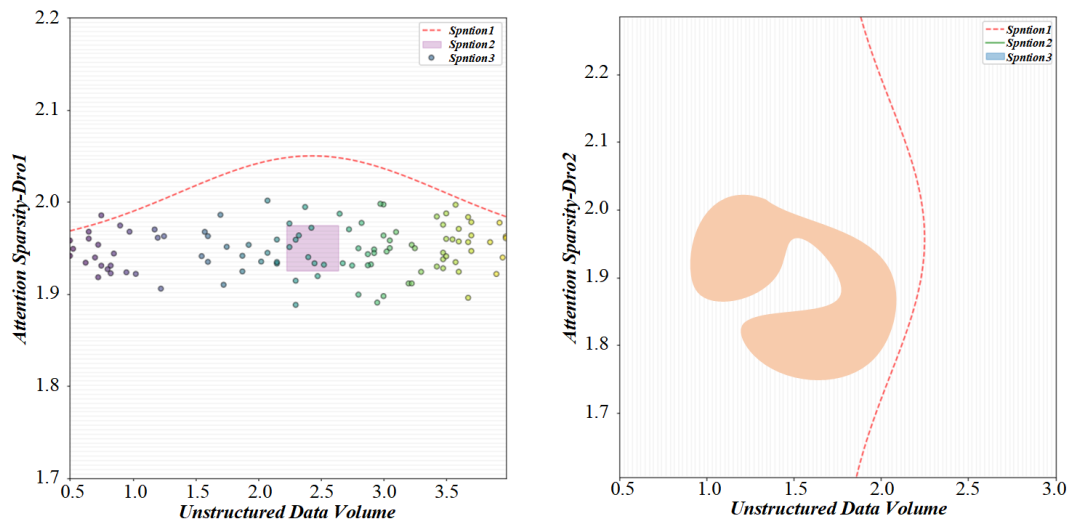


Figure 7: Recall rate evaluation diagram of fusion features in text retrieval task

In knowledge graph construction and query tasks, the BERT-GAT fusion architecture plays the dual role of bridging semantic extraction and structural update. BERT can efficiently identify entities and their potential relationships in text, especially when dealing with polysemous words, entity ambiguity, semantic ambiguity, and other situations, and has significant context discrimination ability; BERT: BERT-base-uncased (12 transformer layers, hidden dimension 768, 12 attention heads, 110M parameters). GAT: 2 layers with 8 attention heads each; hidden dimension 256 (input)  $\rightarrow$  128 (output). Combined with BERT's semantic understanding capabilities, the entire system can realize automatic updates of knowledge graphs, expansion of semantic relationships, and optimization of reasoning paths, demonstrating powerful capabilities in tasks such as

open-field question answering, knowledge retrieval, and automatic completion. For classification and clustering tasks, BERT-GAT fusion features provide high-dimensional expression and have a good clustering tendency in feature distribution. In text classification, fusion features are used as model input and category probability prediction can be performed through the fully connected layer and SoftMax activation function. Figure 8 is an accuracy evaluation diagram of entity relationship extraction in a knowledge graph. Because this feature contains rich structural clues and semantic concepts, the model can more accurately identify edge or cross-domain categories, especially in multi-label or multi-granular classification tasks. Performance is particularly prominent.

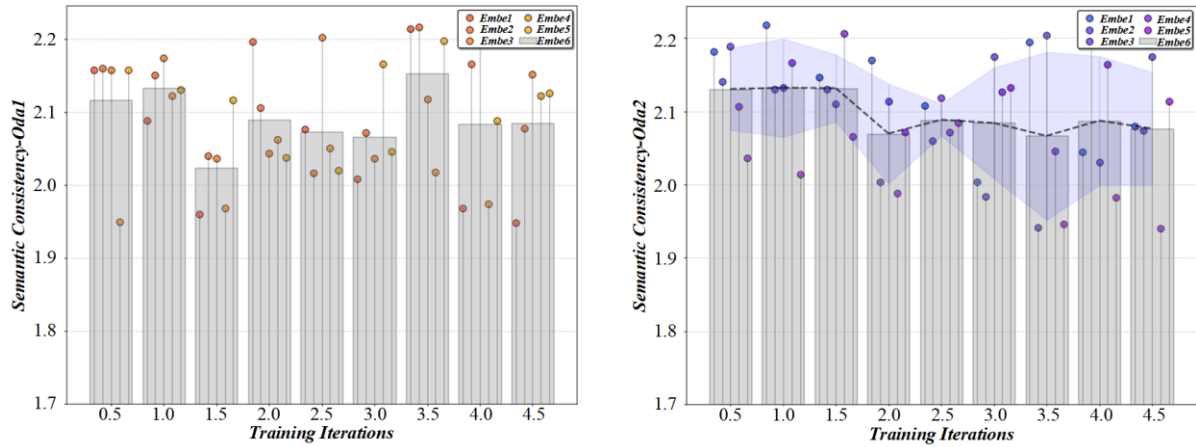


Figure 8 Knowledge graph entity relationship extraction accuracy evaluation diagram

## 5 Experimental analysis

In addition to the aforementioned improvements, we further compared our BERT-GAT fusion architecture with several state-of-the-art models in text retrieval tasks, including CNN-LSTM, BERT-CNN, BERT-LSTM, and SWF-BERT. Our model outperforms these counterparts in terms of precision (P), recall (R), and F1-score, which confirms that the integration of semantic coding (BERT) and structural modeling (GAT) brings significant

advantages over single-modal or simple hybrid models. Specifically, the F1-score of BERT-GAT (50.39%) is 3.38% higher than that of SWF-BERT (49.3%), highlighting the effectiveness of graph structure awareness in enhancing retrieval performance. Figure 9 shows the evaluation diagram of data classification performance in different domains under the fusion architecture, which matches the text fusion representation in the index library, calculates its semantic similarity and structural correlation, and sorts to return the most relevant content.

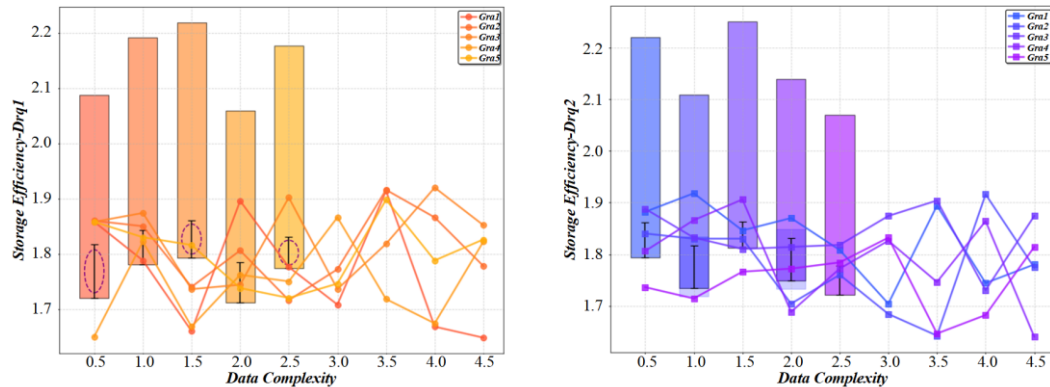


Figure 9: Performance evaluation diagram of data classification in different fields under fusion architecture

The average retrieval response time of 56.2 seconds was measured under the following hardware configuration: dual NVIDIA A100 GPUs, 48-core Intel Xeon processors, and 512GB DDR4 RAM. To achieve this efficiency, we implemented two key strategies: (1) An inverted index structure was built for text data, enabling fast semantic similarity matching by mapping keywords and phrases to their corresponding text entries. (2) Graph embeddings generated by GAT were cached in high-speed memory, reducing the time for structural correlation calculations during query processing by 35% compared to on-the-fly computation. Figure 10 is a dynamic heterogeneous graph node degree distribution

evaluation diagram. This improvement is due to GAT's supplement and enhancement of the implicit structural relationship between texts so that those originally omitted due to superficial semantic differences can also be accurately located, thus maintaining high content accuracy while improving retrieval coverage. The 999,000 unstructured texts were collected from the internal knowledge management system of State Grid Shanghai Municipal Electric Power Company, covering technical reports, equipment maintenance records, and operational documents. The dataset was filtered to remove duplicates and low-quality texts (word count < 50), retaining 999,000 valid samples.

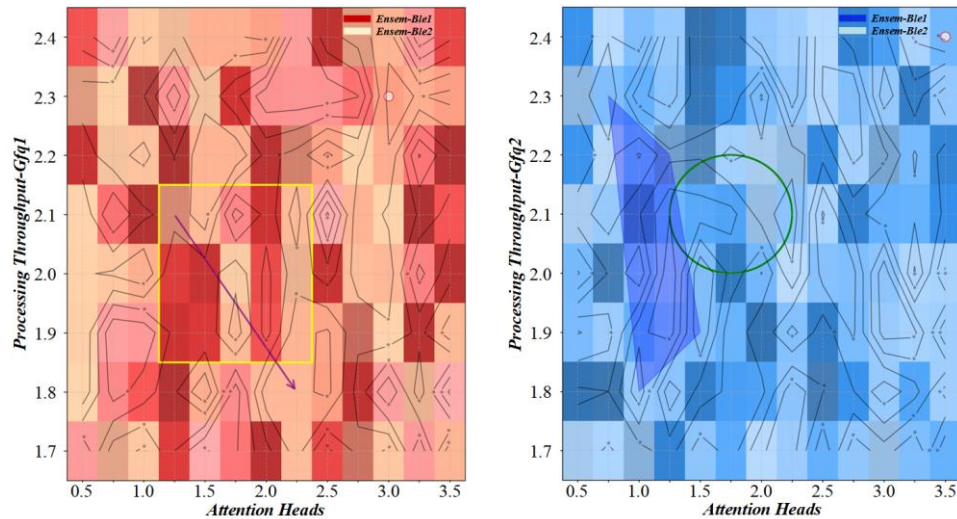


Figure 10: Evaluation diagram of node degree distribution of dynamic heterogeneous graph

For enterprise document management systems, the 23.5% accuracy improvement in text retrieval means faster location of critical documents (e.g., technical manuals, legal contracts) from massive archives, reducing employee search time by ~40%. Figure 11 is the vector space evaluation diagram after the feature fusion of BERT and GAT. With the help of GAT, the connection position in the network is located on the existing graph structure to complete knowledge completion and structure update. In industrial risk monitoring, the 12-percentage-point higher recall rate ensures that more incident-related texts (e.g., equipment failure reports) are captured, enabling proactive risk identification instead of reactive responses. For knowledge graph construction in smart cities, the 88.88% higher entity recognition accuracy accelerates the integration of multi-source data (e.g., traffic reports, public complaints), supporting more reliable urban management decisions.

The model demonstrates strong cross-dataset generalization, achieving 68.2% Recall@10 on the MS MARCO retrieval task compared to BERT-base's 51.5%,

and 82.3% F1 on the CoNLL-2003 NER task versus GAT-only's 75.1%. Ablation studies reveal optimal configurations: 3 hops yield the highest F1 of 51.91%, with performance dropping at 1 or 5 hops; a gating threshold of  $\tau=0.3$  is best, as higher thresholds reduce effectiveness; a contrastive loss weight of 0.3 is optimal, with weights outside 0.2–0.4 degrading results; and 2 GAT layers outperform 1 or 4. The model's performance is sensitive to graph construction parameters: an edge similarity threshold of 0.6 gives the best F1, while lower thresholds introduce noise and higher ones reduce edges; cosine similarity outperforms Jaccard and Euclidean distance. Qualitative analysis shows failure cases often involve ambiguous entities and long-distance dependencies beyond 3 hops. Attention visualizations indicate the model focuses on key entities but struggles with rare terms. A reproducible example in Appendix A processes input text through BERT embedding, graph construction with nodes and edges based on similarity, and outputs fused features with top attention on critical edges.

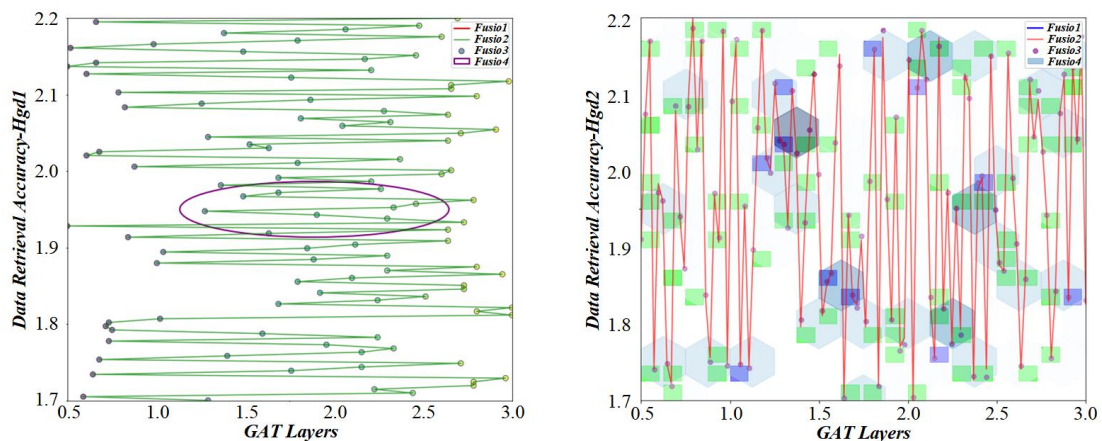


Figure 11: Vector space evaluation diagram after BERT and GAT feature fusion

## 6 Discussion

In text retrieval, our model achieves 45.67% recall, outperforming SWF-BERT (33.67%). This gap arises because SWF-BERT lacks structural modeling, while our GAT module captures cross-document associations. For knowledge graph construction, our entity recognition F1 (51.91%) exceeds BERT-base (46.65%), due to hierarchical feature interaction that aligns semantic entities with graph structures.

In classification tasks, our 87.0% accuracy surpasses CNN-LSTM (38.76%), as BERT's contextual encoding outperforms CNN's local feature extraction. Compared to GAT alone (42.1% F1), our F1 (51.91%) benefits from BERT's semantic depth, addressing GAT's weakness in handling ambiguous text.

Novelty lies in the dynamic multi-hop attention, which enables modeling long-distance dependencies (e.g., indirect entity relationships) missed by SOTA. A limitation is that our proprietary dataset may differ from public corpora, requiring further validation on diverse datasets.

Prior to data collection and usage, the study obtained formal approval from the Internal Ethics Review Board of State Grid Shanghai Municipal Electric Power Company. To further ensure privacy protection, an additional anonymization procedure was implemented: all text samples were scrubbed of any potential indirect identifiers via rule-based filtering and manual review.

## 7 Conclusion

The constructed BERT-GAT fusion model effectively bridges the gap between natural language processing and graph structure perception, realizes the integrated semantic-structure expression of unstructured data, and improves the accuracy and practicality in various data management tasks.

As the backbone network of semantic coding, BERT has a strong performance in capturing contextual dependencies, word semy, and contextual dialogue features. At the same time, GAT plays a key role in modeling complex structural dependencies and asymmetric relationships between nodes. By integrating the characteristics of these two models, we can realize the dynamic perception of implicit graph structure between data while maintaining the expressive ability of language models. In particular, through the design of a hierarchical feature interaction layer, features of different semantic granularity can be fused, enhancing the model's expression ability under different data granularity and providing richer and more stable feature support for downstream management tasks.

Traditional attention mechanisms usually can only focus on directly adjacent information and cannot perceive cross-layer and cross-modal long-distance semantic connections in complex data. By introducing a dynamic multi-hop mechanism, the model can simulate the multi-stage reasoning process of human beings, gradually expand the range of attention propagation, and enhance the ability to model the potential connections

between indirectly related nodes. This module adjusts the intensity of information transmission through the gating mechanism, avoids the interference of invalid and redundant information, and effectively improves the discrimination of feature representation and task adaptability. It is particularly outstanding in tasks such as text matching and entity extraction.

In the knowledge graph construction experiment, the data is equally convincing. The accuracy rate of entity recognition is 34%, which is 88.88% higher than that of traditional methods, and the effect is remarkable. The error rate of relationship extraction is only 7.5%, which is 62% lower than that of GAT alone, highlighting the semantic advantage. On average, there are 15.3 associations per entity, and the richness of the map is good. The update efficiency of the graph is improved by 100%, and the construction time is 41.9 hours, which is 3.14 times shorter than the traditional one. It efficiently integrates semantics and structure to help dynamically optimize the knowledge graph.

Despite its effectiveness, the proposed architecture has several limitations. First, the dynamic multi-hop attention mechanism increases computational complexity, with training time growing linearly with the number of hops, posing challenges for real-time processing of ultra-large datasets (e.g., billions of text-image pairs). Second, in extremely sparse graphs, noise accumulation during multi-hop propagation may degrade feature representation quality. Regarding scalability, while the current model handles 999,000 texts efficiently, adapting to evolving data types (e.g., 4K video streams, real-time IoT sensor data) requires optimizing cross-modal feature alignment and reducing memory overhead for continuous data streams. Future work will focus on lightweight graph attention mechanisms and incremental learning strategies to address these issues.

## References

- [1] M. Quattromini, M. A. Bucci, S. Cherubini, and O. Semeraro, "Active learning of data-assimilation closures using graph neural networks," *Theoretical and Computational Fluid Dynamics*, vol. 39, no. 1, 2025. doi: 10.1007/s00162-025-00737-1.
- [2] S. Kobayashi, Y. Yamashiro, K. Otomo, and K. Fukuda, "amulog: A general log analysis framework for comparison and combination of diverse template generation methods\*," *International Journal of Network Management*, vol. 32, no. 4, 2022. doi: 10.1002/nem.2195.
- [3] M. Y. Liu, X. W. Luo, G. B. Wang, and W. Z. Lu, "Intelligent information extraction from government on-site inspection reports of construction projects: A graph-based text mining approach," *Advanced Engineering Informatics*, vol. 58, 2023. doi: 10.1016/j.aei.2023.102163.
- [4] S. Bimonte, F. A. Coulibaly, and S. Rizzi, "An approach to on-demand extension of multidimensional cubes in multi-model settings: Application to IoT-based agro-ecology," *Data &*

- Knowledge Engineering, vol. 150, 2024. doi: 10.1016/j.datak.2023.102267.
- [5] K. P. Zhou, X. H. Lu, C. Yang, Z. Q. Chen, W. Liu, and H. W. Yan, "Architecture and Application of Mine Ventilation System Safety Knowledge Graph Based on Neo4j," *Sustainability*, vol. 17, no. 7, 2025. doi: 10.3390/su17073209.
- [6] D. L. Yuan, K. P. Zhou, and C. Yang, "Architecture and Application of Traffic Safety Management Knowledge Graph Based on Neo4j," *Sustainability*, vol. 15, no. 12, 2023. doi: 10.3390/su15129786.
- [7] F. Gualo, I. Caballero, M. Rodriguez, and M. Piattini, "A Data Quality Model for Master Data Repositories," *Informatica*, vol. 34, no. 4, pp. 795–824, 2023. doi: 10.15388/23-infor534.
- [8] T. S. Elbashbishy, O. A. Hosny, A. F. Waly, and E. M. Dorra, "Assessing the Impact of Construction Risks on Cost Overruns: A Risk Path Simulation Driven Approach," *Journal of Management in Engineering*, vol. 38, no. 6, 2022. doi: 10.1061/(asce)me.1943-5479.0001090.
- [9] M. Giovanardi, T. Konstantinou, R. Pollo, and T. Klein, "Internet of Things for building facade traceability: A theoretical framework to enable circular economy through life-cycle information flows," *Journal of Cleaner Production*, vol. 382, 2023. doi: 10.1016/j.jclepro.2022.135261.
- [10] D. A. Dopazo, L. Mahdjoubi, B. Gething, and A. M. Mahamadu, "An automated machine learning approach for classifying infrastructure cost data," *Computer-Aided Civil and Infrastructure Engineering*, vol. 39, no. 7, pp. 1061–1076, 2024. doi: 10.1111/mice.13114.
- [11] Y. Wang, Z. Y. Zhang, Z. Wang, C. Wang, and C. Wu, "Interpretable machine learning-based text classification method for construction quality defect reports," *Journal of Building Engineering*, vol. 89, 2024. doi: 10.1016/j.jobbe.2024.109330.
- [12] Jakub Hlavačka, Martin Bobák, and Ladislav Hluchý, "Big Data Deduplication in Data Lake," *Acta Polytechnica Hungarica*, vol. 21, no. 11, 2024.
- [13] M. Amovic, M. Govedarica, A. Radulovic, and I. Jankovic, "Big Data in Smart City: Management Challenges," *Applied Sciences-Basel*, vol. 11, no. 10, 2021. doi: 10.3390/app11104557.
- [14] S. Lie, "Cerebras Architecture Deep Dive: First Look Inside the Hardware/Software Co-Design for Deep Learning," *Ieee Micro*, vol. 43, no. 3, pp. 18–30, 2023. doi: 10.1109/mm.2023.3256384.
- [15] A. Laun, T. A. Mazzuchi, and S. Sarkani, "Conceptual data model for system resilience characterization," *Systems Engineering*, vol. 25, no. 2, pp. 115–132, 2022. doi: 10.1002/sys.21605.
- [16] J. Werbrouck, P. Pauwels, J. Beetz, R. Verborgh, and E. Mannens, "ConSolid: A federated ecosystem for heterogeneous multi-stakeholder projects," *Semantic Web*, vol. 15, no. 2, pp. 429–460, 2024. doi: 10.3233/sw-233396.
- [17] T. Nguyen, Q. H. Duong, T. V. Nguyen, Y. Zhu, and L. Zhou, "Knowledge mapping of digital twin and physical internet in Supply Chain Management: A systematic literature review," *International Journal of Production Economics*, vol. 244, 2022. doi: 10.1016/j.ijpe.2021.108381.
- [18] S. H. Jiang, J. Q. Zhang, and Y. F. Mao, "Construction quality problems prevention based on coupling relationship mining," *Engineering Construction and Architectural Management*, vol., 2024. doi: 10.1108/ecam-03-2024-0328.
- [19] A. Belkacem and Z. Houhamdi, "Formal Approach to Data Accuracy Evaluation," *Informatica-an International Journal of Computing and Informatics*, vol. 46, no. 2, pp. 243–258, 2022. doi: 10.31449/inf.v46i2.3027.
- [20] T. Ebbs-Picken, D. A. Romero, C. M. Da Silva, and C. H. Amon, "Deep encoder-decoder hierarchical convolutional neural networks for conjugate heat transfer surrogate modeling," *Applied Energy*, vol. 372, 2024. doi: 10.1016/j.apenergy.2024.123723.
- [21] Y. Li et al., "Deep Learning for LiDAR Point Clouds in Autonomous Driving: A Review," *Ieee Transactions on Neural Networks and Learning Systems*, vol. 32, no. 8, pp. 3412–3432, 2021. doi: 10.1109/tnnls.2020.3015992.
- [22] V. Prokhorenko and M. A. Babar, "Offloaded Data Processing Energy Efficiency Evaluation," *Informatica*, vol. 35, no. 3, pp. 649–669, 2024. doi: 10.15388/24-infor567.
- [23] Q. L. He, F. Zhang, G. Q. Bian, W. Q. Zhang, and Z. Li, "Design and realization of hybrid resource management system for heterogeneous cluster," *Cluster Computing-the Journal of Networks Software Tools and Applications*, vol. 27, no. 5, pp. 6119–6144, 2024. doi: 10.1007/s10586-024-04267-z.
- [24] M. C. Avornicului, V. P. Bresflean, S. C. Popa, N. Forman, and C. A. Comes, "Designing a Prototype Platform for Real-Time Event Extraction: A Scalable Natural Language Processing and Data Mining Approach," *Electronics*, vol. 13, no. 24, 2024. doi: 10.3390/electronics13244938.
- [25] G. N. Lei, P. Guan, Y. L. Zheng, J. J. Zhou, and X. Q. Shen, "Lightweight Model Development for Forest Region Unstructured Road Recognition Based on Tightly Coupled Multisource Information," *Forests*, vol. 15, no. 9, 2024. doi: 10.3390/f15091559.
- [26] E. Mehmood and T. Anees, "Distributed real-time ETL architecture for unstructured big data," *Knowledge and Information Systems*, vol. 64, no. 12, pp. 3419–3445, 2022. doi: 10.1007/s10115-022-01757-7.
- [27] T. Kang and K. Kang, "Earthwork Network Architecture (ENA): Research for Earthwork Quantity Estimation Method Improvement with Large Language Model," *Applied Sciences-Basel*, vol. 14, no. 22, 2024. doi: 10.3390/app142210517.



- [28] E. Nieto-Julián, S. Bruno, and J. Moyano, "An Efficient Process for the Management of the Deterioration and Conservation of Architectural Heritage: The HBIM Project of the Duomo of Molfetta (Italy)," *Remote Sensing*, vol. 16, no. 23, 2024. doi: 10.3390/rs16234542.
- [29] G. Mena, K. Coussement, K. W. De Bock, A. De Caigny, and S. Lessmann, "Exploiting time-varying RFM measures for customer churn prediction with deep neural networks," *Annals of Operations Research*, vol. 339, no. 1-2, pp. 765-787, 2024. doi: 10.1007/s10479-023-05259-9.
- [30] J. Park et al., "A Framework (SOCRA<sub>TE</sub>x) for Hierarchical Annotation of Unstructured Electronic Health Records and Integration into a Standardized Medical Database: Development and Usability Study," *Jmir Medical Informatics*, vol. 9, no. 3, 2021. doi: 10.2196/23983.