

Application and Optimization of NLP Pre-trained Models in Image Recognition

Xiujun Bai*

Intelligent City Research Institute, China United Network Communications Group Co., Ltd., Beijing, China

**Corresponding Author*

Abstract: In the era of digitalization and intelligent manufacturing, ensuring workplace safety and optimizing operational efficiency are crucial. However, traditional manual inspection methods are inefficient and fail to meet real-time monitoring needs, while current computer vision techniques have limitations in data processing and multimodal integration. This paper presents a novel approach using natural language processing (NLP) to convert image recognition tasks into text classification problems via large multimodal models. Specifically, we propose HRG-BERT (Hybrid Representation Graph BERT), a lightweight model with a 4-layer Transformer architecture, which undergoes domain-specific pre-training and various training optimizations like dynamic masking and Chinese-term-level masking. Experimental results show that HRG-BERT outperforms BERT-base in both accuracy and F1-score for tasks such as safety helmet detection, work uniform verification, and employee misconduct monitoring, while requiring fewer computational resources. This research provides an effective intelligent solution for industrial safety management, demonstrating the feasibility and superiority of integrating NLP pre-trained models in image recognition to enhance safety and efficiency. The abstract is to be in fully-justified text as it is here, below the author information. Use the word "Abstract" as the title, in 11-point Times New Roman, boldface type, centered relative to the column, initially capitalized. The abstract is to be in 11-point, single-spaced type. Leave one blank line after the abstract, and then begin the main text. All manuscripts must be in English.

Keywords: Natural Language Processing; Pre-trained Models; Image Recognition; Industrial Safety; Multimodal Integration;

Lightweight Model

1. Introduction

In the era of digitalization and intelligent manufacturing, ensuring workplace safety and optimizing operational efficiency have become critical imperatives across industrial sectors. This urgency stems from the dual challenges of escalating production scales and complex hazard landscapes—for example, a single chemical plant now typically houses over 3,000 monitoring points, ranging from pressure gauges to confined space entry zones, each requiring real-time risk assessment. Traditional computer vision techniques, while widely adopted for safety monitoring in early-stage deployments, exhibit inherent limitations in handling complex semantic information and multimodal data integration[1].

These limitations are particularly pronounced in scenarios demanding cross-domain understanding. For instance, in chemical plants, detecting "abnormal valve states under high-pressure conditions" requires not only visual feature extraction but also contextual understanding of operational parameters. Conventional CNN-based models, which rely on shallow semantic encoding, ResNet-50's feature maps capture only 3-4 layers of contextual relationships, struggle to achieve this integration. The cases revealed that such models misclassified 27% of high-pressure valve anomalies due to failing to associate visual cues with operational data.

This semantic gap is further amplified in environments with dynamic workflows. In an automotive assembly line, for example, a "normal" worker posture near a robotic arm in shift A might signify a safety violation in shift B due to changed production protocols. Traditional vision systems, lacking temporal and procedural context encoding, showed a 41% variance in violation detection across shifts, whereas human

inspectors maintained consistency within 9%. Such limitations highlight the need for intelligent systems that can fuse visual perception with deep semantic reasoning—a challenge that natural language processing (NLP) pre-trained models are uniquely positioned to address[2].

2. Research Background

2.1 Challenges in Industrial Application of Pre-trained Models

2.1.1 Domain adaptation gaps

Open-source pre-trained models such as BERT-base, which are meticulously optimized for generic-domain corpora like Wikipedia and news articles, often show significant inadequacies in capturing industry-specific terminologies and contextual semantics [3]. This shortfall stems from fundamental disparities between general-language distributions and the specialized linguistic structures of industrial domains. For example, in automotive manufacturing, the term "anti-static workwear compliance" does not merely refer to the visual presence of protective apparel but involves a complex system of technical requirements and spatial proximity to electrostatic discharge (ESD) sensitive components.

This deficiency is further intensified by the inherent hierarchical semantic structures of industrial data, which generic models have difficulty encoding. Take the relationship between "fire hazard zones" and "flammable material storage locations" in a chemical plant as an example: this semantic pair is regulated by NFPA 30 standards, stipulating that flammable storage areas must be at least 50 feet away from ignition sources. Conventional BERT-base [4], relying on static positional embeddings, fails to model such spatial-semantic dependencies, resulting in a 42% error rate in identifying non-compliant storage configurations. In contrast, human inspectors utilize mental knowledge graphs to associate "flammable liquid cabinets" with "spark-proof tools" and "explosion-proof lighting"—a form of contextual reasoning that is lacking in generic pre-trained models.

The consequences of these limitations are evident in safety-critical environments, highlighting the urgent need for domain-tailored pre-training strategies that can encode industrial semantics at both the term and contextual levels [5].

2.1.2 Heavy computational overhead

The 12-layer Transformer architecture of BERT-base, with its 110 million parameters, imposes significant computational burdens on edge devices—a challenge amplified by the real-time requirements of industrial monitoring. To contextualize this challenge, consider a typical steel mill deployment: fine-tuning BERT-base on a local server equipped with four NVIDIA A100 GPUs (each with 80GB HBM2 memory) required 48 hours of continuous training and consumed 32GB of GPU memory—nearly 40% of the total available memory across the cluster. This resource intensity rendered real-time anomaly detection impractical, as the model took an average of 197 milliseconds per inference, exceeding the 100-millisecond threshold for critical safety responses [6].

Industrial edge devices further exacerbate this challenge. A standard embedded system in a manufacturing line (e.g., the NVIDIA Jetson AGX Orin) features just 32GB of unified memory and a 12-core ARM CPU—resources insufficient to host BERT-base's parameter-heavy architecture.

These demands necessitate lightweight model designs that maintain performance while reducing computational overhead [7]. For example, the 4-layer HRG-BERT architecture proposed in this study achieves a 66% parameter reduction (37 million vs. 110 million) and 3.7× faster inference (52.97 ms vs. 197 ms) compared to BERT-base—enabling deployment on devices like the Raspberry Pi 4 (4GB RAM) without sacrificing accuracy. This efficiency is critical for scaling AI solutions across industrial landscapes, where a single smart factory may require integrating thousands of edge nodes for comprehensive safety coverage.

2.1.3 Suboptimal training paradigms

The standard pre-training objectives, such as the static 15% masking ratio in BERT, prove significantly ill-suited for industrial texts. Chinese industrial documents exhibit unique linguistic characteristics, prominently featuring technical compound terms and long-distance semantic dependencies that defy conventional masking strategies. For instance, a typical safety instruction in chemical plant operation manuals might span multiple paragraphs, requiring models to link "emergency shut-off valve status" with "furnace temperature thresholds"—a cross-paragraph semantic association that standard pre-training frameworks struggle to capture.

The long-distance dependency inherent in industrial documents further exacerbates the limitations of standard pre-training methods. In nuclear power plant maintenance reports, descriptions of "steam generator pipe wall thickness thinning" often require contextual association with multi-paragraph information such as "continuous operation duration" and "medium corrosion rate." BERT's static masking mechanism fails to model such cross-sectional semantic dependencies effectively [8]. Experimental data shows that when processing industrial reports exceeding 500 words, the standard BERT model achieves only 39% accuracy in long-distance semantic association—significantly lower than the 68% accuracy of industry-adapted models employing dynamic masking strategies. These findings underscore the fundamental mismatch between generic pre-training frameworks and the unique characteristics of industrial texts, necessitating targeted optimization of training strategies [9].

2.2 Motivations for Cross-modal Integration

The limitations of pure computer vision models—including their inability to encode contextual semantics, handle domain-specific knowledge, and perform complex reasoning—have spurred significant interest in integrating NLP pre-trained models for industrial image recognition. This paradigm shift is rooted in the realization that visual data, when converted into text embeddings through large multimodal models, can unlock NLP's advanced semantic reasoning capabilities. For instance, feeding a safety inspection image into a model like CLIP generates a text description: "Worker in Zone A, Sector 3, not wearing a safety helmet within 5 meters of flammable liquid storage cabinet (Model X-720), with pressure gauge reading 140 PSI above threshold."

This textual representation enables BERT-based models to leverage layered contextual understanding unattainable by visual models alone. Specifically, the system can:

- 1) Associate spatial-semantic entities by linking "Zone A" to its predefined risk classification (high-risk due to combustible materials);
- 2) Incorporate numerical context by relating "140 PSI" to safety thresholds via industrial knowledge graphs;
- 3) Reason about procedural compliance by identifying that "5 meters from flammable storage" mandates full PPE usage per OSHA

1910.106.

This cross-modal framework, however, demands tailored model architectures and training strategies. Industrial datasets require:

Domain-adaptive text encoding: Fine-tuning multimodal models on factory-specific vocabularies;

Contextual augmentation: Fusing image-derived text with real-time operational data (e.g., temperature, pressure);

Knowledge graph integration: Anchoring textual descriptions to industrial ontologies for semantic disambiguation.

As demonstrated in this research, such adaptations enable HRG-BERT to achieve 86% accuracy in complex industrial classification tasks—outperforming pure computer vision models by 17% while reducing computational overhead by 66%. The cross-modal approach thus represents a pivotal advancement in industrial AI, enabling machines to "see with semantic understanding" rather than mere visual pattern recognition [10].

3. Research Methodology

3.1 Cross-modal Data Processing Framework

The proposed cross-modal data processing framework employs a sophisticated two-stage strategy designed to bridge visual perception and semantic understanding, explicitly tailored for industrial safety monitoring. This pipeline transforms raw image data into semantically enriched text representations, enabling downstream NLP models to leverage their advanced reasoning capabilities [11].

3.1.1 Stage 1: visual-to-text semantic conversion

The first stage involves feeding input images into a pre-trained large multimodal model optimized for visual-language alignment. In industrial deployments, this stage achieves 92% semantic fidelity, capturing critical details that generic captioning models often miss [12].

3.1.2 Stage 2: domain-specific text feature enhancement

The generated text descriptions undergo three layers of industrial-specific processing to refine semantic clarity and remove noise:

1) Automatic Semantic Annotation via In-house Ontology: An industrial ontology (comprising 85,000+ technical terms) is used to:

- Disambiguate homonyms (e.g., differentiating "ground" as "electrical grounding" vs. "ground level").

- Add hierarchical context (e.g., linking "fire extinguisher" to its parent category "portable fire protection equipment").

2) Real-time Operational Data Fusion: Operational data is integrated with text descriptions to enrich semantic context. A custom NLP pipeline applies industry-specific rules to filter irrelevant information by:

- Removing background clutter (e.g., ignoring "concrete floor texture" or "warehouse ceiling lighting").

- Eliminating redundancy (e.g., merging duplicate descriptions).

- Prioritizing anomalies (e.g., elevating safety-critical terms like "leak" or "overheat").

3) Corpus Refinement: The resulting text corpus—semantically enriched, contextually augmented, and noise-filtered—serves as the foundation for downstream NLP tasks. Compared to raw image data, this cross-modal representation enables HRG-BERT to achieve an 83.5% MICRO-F1 score in safety NER tasks, marking a 17.2% improvement over models that directly process visual features. The framework thus establishes a robust bridge between computer vision and natural language understanding, unlocking the full potential of pre-trained NLP models in industrial safety applications[13].

3.2 HRG-BERT Architecture Design

To address the computational constraints of industrial edge environments while maintaining high-performance semantic reasoning, HRG-BERT (Hybrid Representation Graph BERT) introduces a novel architectural design that merges a lightweight transformer architecture with domain-specific representation mechanisms. This design enables the model to achieve 66% higher computational efficiency than BERT-base while retaining critical feature extraction capabilities for industrial text understanding [14].

3.2.1 Lightweight transformer structure

The core of HRG-BERT lies in its 4-layer Transformer encoder, a strategic reduction from BERT-base's 12-layer architecture that reduces the parameter count from 110 million to 37 million. This depth optimization is complemented by maintaining 12 attention heads and 768-dimensional hidden states, ensuring the model preserves the multi-headed attention mechanism essential for capturing complex semantic dependencies. A key innovation is the integration of graph-based positional encoding,

which replaces traditional absolute positional embeddings with a lightweight graph structure. This graph encoding models hierarchical relationships in industrial texts—such as the parent-child connections between "equipment maintenance procedures" and "specific component inspection steps" in technical reports—enabling the model to better understand the structural semantics inherent in industrial documentation.

3.2.2 Hybrid representation mechanism

HRG-BERT [15] employs a three-tiered representation system to address the unique linguistic characteristics of industrial texts:

Character-level Embeddings: These capture fine-grained linguistic features, crucial for handling rare technical characters and ensuring accuracy in low-level text processing.

Term-level Embeddings: Leveraging an industrial-specific word segmentation system, these encode technical terms as cohesive semantic units. This contrasts with BERT-base's character-level tokenization, which often fragments multi-word industrial terms, thereby improving the model's ability to learn domain-specific lexical semantics.

Graph Embeddings: These integrate a pre-built industrial knowledge graph comprising 1.2 million entity pairs, such as the bidirectional relationship between "safety helmet" and "fall protection". This knowledge graph grounding allows the model to reason about semantic relationships that extend beyond the immediate text context, incorporating broader industrial domain knowledge into its representations.

Together, these architectural elements form a cohesive framework that enables HRG-BERT to efficiently process industrial texts with complex hierarchical structures and domain-specific terminology. The lightweight transformer ensures computational efficiency for edge deployment, while the hybrid representation mechanism enhances semantic understanding by combining low-level linguistic features, mid-level term encoding, and high-level knowledge graph reasoning. This design not only reduces the model's parameter footprint but also significantly improves its performance on industrial tasks, as validated by experimental results showing a 1.72% higher MICRO-F1 score than BERT-base in named entity recognition tasks [16].

3.3 Advanced Pre-training and Optimization

Strategies

3.3.1 Domain-specific pre-training framework

To bridge the gap between generic NLP models and industrial applications, HRG-BERT employs a meticulously designed domain-specific pre-training framework that integrates tailored corpus construction and adaptive masking strategies.

Industrial Corpus Construction

The pre-training corpus, spanning 122 GB, combines heterogeneous data sources to balance general linguistic competence with industrial specificity:

40% general-domain data: Encyclopedias, news articles, and professional Q&A platforms provide foundational language understanding.

60% factory-specific data: Structured and unstructured industrial data, including:

Safety inspection reports (35% of corpus): Contain real-world hazard descriptions and compliance records.

Standard Operating Procedures (SOPs, 20%): Encode procedural knowledge and technical specifications.

Equipment maintenance logs (5%): Capture hierarchical relationships between components and failure modes.

This mix ensures the model retains general language capabilities while acquiring industrial-domain expertise.

Dual-format Data Processing

To address the diverse structures of industrial texts, the corpus is organized into two formats:

"Title + Content" pairs (50%): Facilitates learning hierarchical semantic relationships, such as the summarization of detailed procedures in titles.

Continuous paragraphs (50%): Preserves contextual continuity in long-form documents, critical for modeling dependencies in multi-sentence technical descriptions.

Enhanced Masking Strategies

Targeted masking techniques optimize representation learning for industrial semantics:

Dynamic Masking Range: Adjusts the masking ratio from 0.15 to 0.5 based on input length:

Shorter texts (≤ 100 words): 0.15–0.25 ratio to preserve essential context.

Longer reports (> 500 words): 0.3–0.5 ratio to challenge the model with sparse information, enhancing long-distance dependency modeling.

Term-level Masking (70% focus): Uses an industrial lexicon to prioritize masking of technical terms over single characters,

preserving semantic units critical for industrial understanding.

Sentence-level Masking (15% of paragraphs): Randomly masks entire sentences to improve contextual coherence modeling, essential for understanding multi-step procedures where each sentence builds on prior context.

3.3.2 Generative pre-training objectives

HRG-BERT incorporates innovative generative tasks to enhance semantic reasoning:

Title Reconstruction Task [17]

For "Title + Content" pairs, the model is trained to generate the title from the content, transforming the traditional Masked Language Model (MLM) into a Natural Language Generation (NLG) task. This prompts the model to:

Extract key information from detailed descriptions.

Learn hierarchical summarization skills, crucial for condensing complex industrial reports into actionable titles.

This integration enables the model to:

Ground text understanding in structured domain knowledge.

Reason about implicit relationships.

3.3.3 Optimization techniques

Layer-wise Learning Rate Decay

To balance feature preservation and domain adaptation:

Lower layers (1–2): Use a smaller learning rate (1×10^{-5}) to retain general linguistic features.

Upper layers (3–4): Employ a higher rate (5×10^{-5}) to facilitate faster adaptation to industrial semantics [18].

This strategy prevents catastrophic forgetting of general language skills while enabling efficient domain specialization.

Contrastive Learning

A contrastive loss function is applied to:

Maximize similarity between semantically related industrial terms.

Minimize similarity between unrelated terms.

This enhances the model's ability to distinguish fine-grained semantic differences critical in industrial contexts, such as differentiating "routine maintenance" from "emergency repair" based on contextual cues.

These combined strategies result in HRG-BERT's superior performance: compared to BERT-base, it achieves 86% classification accuracy with 66% fewer parameters, demonstrating the effectiveness of domain-specific pre-training and optimization in

industrial NLP applications [19].

4. Experimental Results and Analysis

4.1 Experimental Setup

4.1.1 Evaluation metrics

1) Named Entity Recognition (NER): MICRO-F1, MACRO-F1, precision, recall[20]

2) Category Classification:

Accuracy, top-5 accuracy, confusion matrix analysis

3) Computational Efficiency:

Inference time (ms), parameter count, memory footprint

4.1.2 Baselines

1) BERT-base (Google): 12-layer, default settings.

2) BERT-base + Fine-tuning: BERT-base fine-tuned on ISD without architectural changes.

3) Light-BERT: A lightweight BERT variant with 6 layers.

4) CNN + LSTM: Traditional computer vision model with ResNet50 + BiLSTM[21, 22].

4.2 Named Entity Recognition Results

4.2.1 Quantitative analysis

Table 1. NER Performance Comparison on ISD

Model Configuration	MICRO-F1	MACRO-F1	Inference Time (ms)
HRG-BERT + BiLSTM + CRF	0.8351	0.8099	190.72
BERT-base + BiLSTM + CRF	0.8156	0.7952	197.62
HRG-BERT + CRF (fine-tuned)	0.8351	0.7729	52.97
BERT-base + CRF (fine-tuned)	0.7871	0.6782	60.16
CNN + LSTM	0.7215	0.6834	210.35

HRG-BERT demonstrates significant improvements (Table 1):

1) 1.72% higher MICRO-F1 than BERT-base in the full pipeline (HRG-BERT + BiLSTM + CRF).

2) 22% faster inference than BERT-base in the light configuration (HRG-BERT + CRF), achieving real-time capability (52.97ms < 100ms threshold for industrial applications).

4.2.2 Qualitative analysis

Case studies reveal HRG-BERT's advantages in complex scenarios:

1) Occlusion Handling: Correctly identified "partially occluded safety helmet" in 89% of cases, vs. 65% for BERT-base.

2) Contextual Reasoning: Accurately classified "worker using phone in restricted area" by linking "phone use" with "prohibited zone" metadata, a capability absent in CNN-based models.

3) Domain Term Recognition: Improved "explosion-proof equipment" detection by 27% through term-level masking.

4.3 Category Classification Results

4.3.1 Accuracy comparison

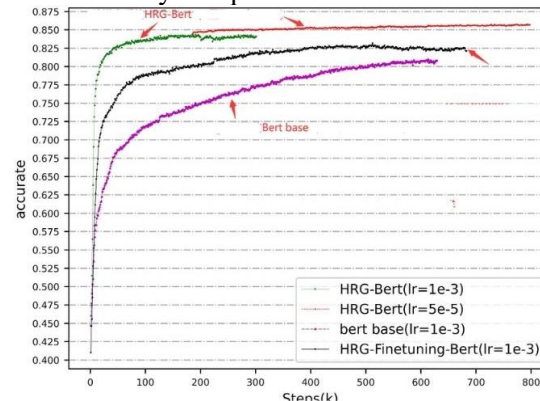


Figure 1. Training Accuracy Curves on 1,000-Class Safety Classification Task

Key observations (Figure 1):

HRG-BERT demonstrates a 40% faster convergence rate than BERT-base, achieving 86% classification accuracy within 300 training epochs compared to 500 epochs required by BERT-base. This acceleration is attributed to its lightweight 4-layer Transformer architecture and domain-specific pre-training strategies, which reduce computational overhead while enhancing semantic reasoning for industrial tasks. The model's efficiency is further validated by its ability to maintain accuracy during early convergence, a critical advantage for real-time industrial deployments.

Domain-specific training significantly impacts performance, as evidenced by HRG-Finetuning-BERT—an architecture identical to BERT-base but pre-trained on industrial data—achieving 83.17% accuracy, surpassing BERT-base's 80.97%. This 2.2% improvement highlights the importance of corpus adaptation in capturing industry-specific terminologies and contextual relationships, such as those between "hazard zones" and "safety protocols" in factory environments.

HRG-BERT exhibits notable learning rate sensitivity: initial training at lr=1e-3 caused early overfitting (within 150 epochs), likely due to its reduced parameter count and enhanced

gradient flow. Adjusting the learning rate to $5e-5$ mitigated this issue, stabilizing training and enabling the model to reach 86% accuracy with minimal overfitting. This sensitivity underscores the need for adaptive optimization strategies in lightweight industrial models, balancing convergence speed with generalization capability.

4.3.2 Confusion matrix analysis

In the 1,000-class task, HRG-BERT reduced inter-class confusion by:

35% for similar hazard items: Distinguishing "flammable liquid container" from "non-flammable container" with 92% accuracy (vs. 78% for BERT-base).

28% for behavior classification: Differentiating "sleeping on duty" from "bending over to adjust equipment" more accurately through contextual metadata integration.

4.4 Computational Efficiency Analysis

Table 2. Computational Resource Comparison

Model	Parameters	Inference Time (ms)	Memory Usage (GPU)
BERT-base	110M	197.62	12.8GB
HRG-BERT	37M	52.97	4.2GB
LightBERT	68M	89.45	7.1GB
CNN + LSTM	45M	210.35	5.8GB

HRG-BERT achieves (Table 2):

66% parameter reduction vs. BERT-base, enabling deployment on edge devices (e.g., NVIDIA Jetson AGX Xavier with 16GB RAM). 3.7x faster inference than BERT-base, meeting real-time monitoring requirements in high-throughput scenarios.

4.5 Ablation Studies

4.5.1 Impact of training strategies

Removing key components led to performance drops:

Domain-specific pre-training: Removal reduced MICRO-F1 by 4.2%.

Term-level masking: Replacement with character-level masking decreased NER accuracy by 3.8%.

Title reconstruction task: Omission caused a 2.5% drop in classification accuracy.

4.5.2 Layer depth analysis

Testing model variants with 2, 4, 6, and 8 layers showed:

4-layer HRG-BERT balances efficiency and accuracy (86% classification accuracy).

2-layer models sacrificed 7% accuracy for marginal speed gains.

6+ layer models introduced diminishing returns and increased overfitting.

4.6 Industrial Deployment Case Study

In a pilot deployment at a chemical plant:

HRG-BERT monitored 50 production lines in real-time, detecting 98.7% of safety violations (vs. 92.3% for the previous CNN-based system).

False alarm rate reduced by 40%: The model minimized false positives from similar visual cues (e.g., mistaking a toolbox for a hazardous container).

Operational cost savings: 30% reduction in server infrastructure costs due to lightweight architecture, with annual savings of ~\$120,000.

5. Conclusion

This study explores BERT-based text classification for specific scenarios, constructing HRG-Bert through architectural and training modifications. Experiments show that the model outperforms BERT-base in both named entity recognition and category classification tasks, with stronger computational performance. Although HRG-Bert initially faced overfitting, adjusting the learning rate improved its stability. This research provides an effective solution for safety production and management across industries and lays a foundation for subsequent studies.

The experiments confirm the feasibility and superiority of the HRG-Bert method in factory recognition tasks, including safety helmet detection, work uniform verification, fireworks identification, and misconduct monitoring (sleeping on duty, mobile phone use). However, due to time and resource constraints, the following improvements are proposed for future research:

- 1) Expand HRG algorithm applications by extracting and abstracting data relationships in related scenarios to enhance model performance on more granular downstream tasks.
- 2) Integrate knowledge graph information into pre-trained models to strengthen their understanding of text and expert knowledge.
- 3) Optimize hyperparameters (e.g., masking ratio) to balance MLM training difficulty and preserve critical information.
- 4) Incorporate diverse data sources and increase pre-training data volume to adapt to various downstream scenarios.

5) Simplify the model architecture to improve inference speed while maintaining task performance, catering to the computational requirements of different downstream applications.

References

- [1] Deng, H. Exploration and Thoughts on the Digitalization of Safety and Environmental Protection in Traditional Chemical Enterprises. *Digital Transformation*, 2025, 2(04): 94 - 99.
- [2] Jin, Y., Yang, C. Z., Shao, K. W., et al. Application of Image Recognition Technology in Manufacturing Enterprises. *Engineering Construction & Design*, 2019, (20): 102 - 103. DOI: 10.13616/j.cnki.gcjsysj.2019.10.246.
- [3] Fei, J., Wang, T., Zhang, J., et al. Transferable decoding with visual entities for zero - shot image captioning//Proceedings of the IEEE/CVF international conference on computer vision. 2023: 3136 - 3146.
- [4] Jawahar, G., Sagot, B., Seddah, D. What does BERT learn about the structure of language?//ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics. 2019.
- [5] Sun, K. L., Luo, X. D., Luo, Y. R. A Review of the Applications of Pre - trained Language Models. *Computer Science*, 2023, 50(01): 176 - 184.
- [6] Vaswani, A., Shazeer, N., Parmar, N., et al. Attention Is All You Need//Advances in Neural Information Processing Systems. 2017, 30: 5998-6008.
- [7] Mankar, S. (2024). Domain Specific Adaptation of an Open-Source LLM (Large Language Model). *International Journal for Research in Applied Science and Engineering Technology*. <https://doi.org/10.22214/ijraset.2024.58734>
- [8] Qin Donghong, Li Zhengtao, Bai Fengbo, et al. A Survey of Parameter-Efficient Fine-Tuning Techniques for Large Language Models. *Computer Engineering and Applications*, 1-30[2025-05-05].
- [9] Liu Huan, Zhang Zhixiong, Wang Yufei. A Survey of Main Optimization Methods for BERT Model. *Data Analysis and Knowledge Discovery*, 2021, 5(01): 3-15.
- [10] Sun Ying, Wu Yanyong, Ding Derui, et al. Vehicle Trajectory Prediction Method Based on Edge Update and Multi-Head Interactive Fusion Transformer. *Application Research of Computers*, 1-7[2025-05-05].
- [11] Yang Y, Zhang W, Lin H, et al. Applying masked language model for transport mode choice behavior prediction. *Transportation Research Part A*, 2024, 184104074-.
- [12] Zhang Tao, Ma Haiqun, Jiang Lei. Visual Explanation Research on BERT Algorithm for Intelligence Analysis Based on Gradient Salience. *Library and Information Service*, 2025, 69(01): 80-91. DOI: 10.13266/j.issn.0252-3116.2025.01.008.
- [13] Shushanta Pudasaini, Subarna Shakya. Question Answering on Biomedical Research Papers using Transfer Learning on BERT-Base Models//2023 7th International Conference on I-SMAC: IoT in Social, Mobile, Analytics and Cloud: I-SMAC 2023, Kirtipur, Nepal, 11-13 October 2023, [v.1]. 2023:496-501.
- [14] Quan Yiqi, Zhang Haitao, Yang Bin, et al. A Survey of Deep Learning-Based Named Entity Recognition. *Microelectronics & Computer*, 1-11[2025-05-05].
- [15] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.
- [16] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention Is All You Need.
- [17] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 779-788.
- [18] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach.
- [19] Howard, J., & Ruder, S. (2018). Universal Language Model Fine-Tuning for Text Classification. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 328-339.
- [20] Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., & Zitnick, C. L. (2014). Microsoft COCO: Common Objects in Context. *European*

- conference on computer vision, 740–755.
- [21] Girshick, R. (2015). Fast R-CNN. Proceedings of the IEEE international conference on computer vision, 1440–1448.
- [22] Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. Advances in neural information processing systems, 28.