# Fine-Tuning the Qwen2.5-VL Model for Intelligent Applications in the Electrical Domain

Yao Song[1, 2,*] and Chunli Lv[2] and Kun Zhu[3] and Xiaobin Qiu[1]

[1] Information Office, China Agricultural University, Beijing, 100083, China
[2] College of Information and Electrical Engineering, China Agricultural University, Beijing, 100083, China
[3] China Petroleum Engineering Construction Co., Ltd. Beijing, 100083, China

## Abstract

This study explores the fine-tuning application of the Qwen2.5-VL multi modal large model in the electrical domain. The electrical industry faces numerous challenges in maintaining and managing complex electrical systems. Traditional methods often rely on manual inspection and analysis. With the rapid advancement of artificial intelligence (AI) technologies, there is a growing need to explore how these tools can be applied to improve efficiency and accuracy in the electrical domain. Qwen2.5-VL is a state-of-the-art visual language model. We adopted the LoRA (Low Rank Adaptive) method to fine tune the model, which enables efficient parameter updates in low resource environments while maintaining high performance. This study analyzes the data characteristics and task requirements in the electrical domain, designs fine-tuning strategies with a focus on image-based applications, including data preprocessing, model fine-tuning, and training parameter optimization. The experimental re-sults show that the fine tuned model has achieved significant performance im-provements in tasks such as electrical equipment fault detection, image recogni-tion, and text classification. This study provides new ideas and methods for the application of artificial intelligence in the electrical domain, which is of great significance for promoting the development of electrical intelligence.

## 1. Introduction

With the rapid development of artificial intelligence technology, the application of large models in various fields is becoming increasingly widespread [1]. Electrical engineering, as one of the core areas of modern industry, is facing challenges and opportunities in digital transformation [2]. Applying advanced artificial intelligence technology to the electrical field is of great significance for improving production efficiency, enhancing equipment reliability, optimizing resource management, and promoting sustainable development [3]. As a new generation multimodal large model, Qwen2.5-VL has powerful language understanding and image processing capabilities, providing new possibilities for intelligent applications in the electrical field [4].

The electrical field is facing challenges such as equipment fault detection, image recognition in complex environments, and multimodal data fusion [5]. The existing computer vision systems are unable to cope with the complexity of electrical equipment, including complex electrical components, dynamic environments, and the need for multimodal data fusion [6]. Therefore, developing specialized models for the electrical field is of great significance [7].

The integration of artificial intelligence (AI) into the electrical domain has revolutionized traditional practices,

*Corresponding author: Email: songyao@cau.edu.cn

enabling advanced applications such as fault diagnosis, power grid optimization, and equipment condition monitoring [8]. Among AI technologies, vision-language models (VLMs), which combine visual and textual understanding, hold significant potential for interpreting complex electrical data, including schematic diagrams, infrared images, and equipment manuals [9]. However, the deployment of general-purpose VLMs in domain-specific scenarios often faces challenges due to the unique characteristics of electrical systems, such as specialized terminology, intricate topological relationships, and safety-critical requirements [10].

Globally, VLMs like GPT-4V and BLIP-2 have demonstrated remarkable performance in cross-modal tasks, including image captioning and visual question answering [11]. In China, models such as Qwen-VL and ERNIE-ViL have also achieved breakthroughs in integrating domain knowledge with multimodal reasoning. Recent efforts to adapt VLMs to specialized fields include medical image analysis and autonomous driving, where domain-specific fine-tuning strategies and knowledge injection have proven effective [12]. In the electrical domain, preliminary studies have explored convolutional neural networks and transformers for tasks like insulator defect detection and load forecasting [13]. However, these works predominantly focus on single-modal data (images or text), neglecting the synergistic analysis of multimodal information inherent in electrical systems [14]. For instance, interpreting a substation's operational status often requires correlating thermal images with maintenance logs—a capability underdeveloped in existing approaches.

Three critical limitations hinder the application of general VLMs in the electrical domain:

(1) Domain-Specific Knowledge Gap: Pretrained VLMs lack familiarity with electrical terminologies ("partial discharge" or "phasor measurement units") and structured data formats (single-line diagrams), leading to suboptimal performance in semantic alignment [15].

(2) Data Complexity and Scarcity: Electrical datas often involve high-resolution images, symbolic notations, and heterogeneous formats, yet publicly available multimodal datasets tailored to this domain are scarce [16, 17].

(3) Inefficient Fine-Tuning Paradigms: Conventional fine-tuning methods, designed for generic scenarios, struggle to preserve the model's generalizability while adapting to specialized tasks, risking catastrophic forgetting or overfitting[18].

This study aims to bridge these gaps by developing a domain adaptive fine-tuning framework for Qwen2.5-VL, which is specifically optimized for electrical applications. The main objectives include:

(1) Building a multimodal electrical dataset: Organize a dataset that includes images (such as device snapshots, infrared thermography), textual descriptions (such as manuals, fault reports), and structured data (such as circuit diagrams) to capture domain specific features.

(2) Resource optimization: Develop efficient fine-tuning strategies to reduce computational resource requirements, making the model more suitable for practical applications in the electrical field, and improving the performance of specific tasks without compromising the basic functionality of the model.

(3) Enhance accuracy: By fine-tuning data in the electrical field, improve the accuracy of the model in electrical equipment fault detection and image recognition tasks.

(4) Domain adaptability: By fine-tuning a specialized model for the electrical field, the model can be seamlessly integrated into existing electrical management systems to improve interpretability in scenarios such as fault location.

(5) Verify actual effectiveness: Evaluate fine-tuning models in practical tasks, including device state recognition, cross pattern retrieval of maintenance records, and security violation detection.

By addressing these challenges, this research seeks to establish a robust framework for deploying VLMs in intelligent electrical systems, advancing both AI methodology and energy infrastructure management. The outcomes are expected to provide insights into domain-specific adaptation of multimodal models, with implications for industrial AI applications beyond the electrical sector.

## 2.    Related Methods

This study uses the Qwen2.5-VL-7B-Instruction model, which is a powerful multimodal large-scale language model used to solve specific tasks in the electrical field. And the LoRA method is used for fine-tuning, which is known for its high efficiency in parameter efficient fine-tuning (PEFT) of large models [19]. The fine-tuning process is based on the organized electrical equipment image dataset, which is a comprehensive collection of images and annotation information of a large number of electrical equipment.

(1) Data preprocessing: The data preprocessing stage mainly collects image and text data of electrical equipment, annotates and preprocesses them, including image standardization and text segmentation processing [20].

(2) Training parameter optimization: Various parameter adjustment strategies are adopted, combined with hybrid training and distributed training techniques, to improve training efficiency.

Qwen2.5-VL-7B-Instruction is a visual language model that performs well in tasks involving text and image data [21]. It performs well in various benchmark tests, including visual Q&A and document comprehension. The reason for choosing this model is that it can handle multimodal inputs and adapt to professional fields such as electrical engineering.

The dataset contains various electrical professional images, annotations, and related metadata [22]. It is carefully planned to address unique challenges in the electrical field, such as electrical equipment identification and environmental analysis. This dataset contains high-quality images and detailed annotations, suitable for fine-tuning large models such as Qwen2.5-VL-7B-Instruction.

The fine-tuning training process is based on LoRA. LoRA is a parameter-efficient fine-tuning method designed for large language models. The LoRA method was used to fine-tune the Qwen2.5-VL-7B-Instruct model. The core idea is to freeze the pre-trained model weights and inject low-rank

matrices to specific layers of the model, allowing for efficient updates without modifying the entire model architecture [23]. The basic principle of LoRA training process is shown in Figure 1.
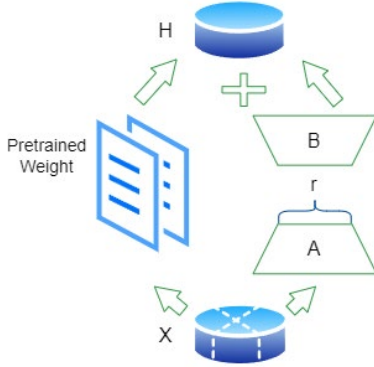


**Figure. 1.** The basic principle of LoRA fine-tuning process

The fine-tuning process only trains A and B. The key steps of the fine-tuning process are as follows:

(1) Model Initialization: The Qwen2.5-VL-7B-Instruct model was loaded and prepared for fine-tuning. The model's weights were frozen except for the layers designated for LoRA adaptation.

(2) Data Preprocessing: The dataset was preprocessed to convert the image and text data into a format compatible with the model. This involved resizing images, tokenizing text, and creating input tensors.

(3) LoRA Configuration: The LoRA configuration was set up to target specific layers of the model, such as the query [24], key, and value projections. The parameters for LoRA, including rank, alpha, and dropout, were carefully tuned to optimize performance.

(4) Training: The model was fine-tuned using the LoRA method with a training configuration that included a learning rate of $1 \times 10{-4}$, a batch size of 4, and gradient accumulation steps to manage computational resources. The training process involved multiple epochs, with periodic evaluation and checkpointing to monitor progress.

(5) Evaluation: After fine-tuning, the model was evaluated on a validation set from the dataset to assess its performance in tasks such as electric equipment identification and classification. The evaluation metrics included accuracy, precision, recall, and F1-score [25].

By leveraging the LoRA method, we were able to efficiently adapt the Qwen2.5-VL-7B-Instruct model to the electric domain, achieving significant improvements in task-specific performance while minimizing computational overhead.

In the data preprocessing stage, we collected and integrated multimodal data in the power field, including equipment technical documents, fault reports, scientific research papers, high-resolution equipment images (such as infrared thermal imaging, circuit topology diagrams), and sensor timing data [26]. For textual data, we have standardized power professional terminology, filtered noise, and annotated fine-

grained entities (including equipment models, fault types, etc.). The image data undergoes standardization processing (unified resolution and color space) and enhancement operations (simulating device states under different lighting conditions, adding noise to enhance robustness). In addition, for power time series data (such as voltage fluctuation records)[27], we use sliding window segmentation and normalization processing to extract spatiotemporal features.

At the model architecture level, we retained the core multimodal alignment capability of Qwen2.5-VL, but made targeted improvements for the characteristics of the power field.

To improve fine-tuning efficiency and model performance, we have designed the following training strategies. Progressive learning rate scheduling, initially using a low learning rate (1e-3) to stably adapt to the distribution of power data, and then gradually increasing to the peak to accelerate convergence and avoid gradient oscillations caused by domain differences. Adopting mixed precision training and distributed parallelism: utilizing mixed precision to reduce video memory usage and support larger batch inputs. Domain adaptive data augmentation is used to address the scarcity of power data. Synthetic data generation is employed, such as simulating equipment images with different levels of faults, and integrating contrastive learning to enhance the robustness of cross modal representations.

The above strategies aim to address the unique challenges in the field of electricity. Professionalism of terminology: Improve semantic parsing accuracy through domain dictionary injection and fine-grained annotation. Data heterogeneity: The spatiotemporal fusion module and multimodal alignment optimization support complex scenarios. Security sensitive requirements: Regularization and synthetic data augmentation ensure the model's generalization ability in limited samples and reduce the risk of misjudgment.

This fine-tuning scheme can significantly improve the performance of Qwen2.5-VL in power tasks, such as equipment status monitoring, combining infrared images with operation logs, cross modal fault retrieval, matching historical fault images with text queries, safety risk prediction, and joint inference of temporal data and visual features.

## 3. Data Analysis

The image data in the electrical field has the characteristics of diversity, complexity, and high annotation requirements. The diversity of data is reflected in the differences in the types of electrical equipment and working environments; complexity is reflected in the complex structure of the device and the dynamic environmental background. Therefore, accurate labeling is crucial for model training.

The image data characteristics are diverse in the field of electricity [28]. The image of power equipment covers multiple dimensions, including equipment types such as transformers, circuit breakers, insulators, normal or faulty operating status, and detection modes such as visible light and infrared thermal imaging [29]. For example, the same device will exhibit significantly different visual characteristics under

different fault modes, such as partial discharge and mechanical deformation.

The image data characteristics are complex in the field of electricity. The appearance of the equipment is significantly affected by environmental factors [30], such as changes in daytime and nighttime illumination, and attenuation of infrared images in rainy and snowy weather. Multiple background interferences: Images often contain complex scenes [31], such as multiple equipment stacks in substations and staggered cable layouts, requiring differentiation between target devices and background elements such as brackets and vegetation. Coexistence of multi-scale features: There is a need to capture both macroscopic overall states [32], such as the overall distribution of equipment, and microscopic defects, such as surface cracks on insulators.

High annotation are required. It is necessary to accurately label the equipment type, fault type such as arc, overheating, mechanical damage, severity level and location information, and defect area boundary [33]. Cross modal alignment is also crucial. Some images need to be associated with textual descriptions, such as fault descriptions in maintenance reports [34].

Construction process of power image data is important. It includes collecting power company operation and maintenance databases, public datasets, shared data from research institutions, and on-site filming. Types include device appearance images, infrared thermal images, discharge detection spectra, circuit schematics, etc [35]. Data annotation process is to annotate the main body of the equipment, such as transformer oil pillows, defect areas, such as insulator cracks, safety signs, such as high voltage warning signs [36]. Attribute annotation includes recording device model and operating parameters, such as load rate, ambient temperature, and fault codes.

Data preprocessing is to perform image processing, unify resolution, and align multispectral image channels [37]. Simulated data augmentation includes adding noise and blurring. Text processing includes professional terminology cleaning,

such as unifying circuit breaker and CB expressions. Entity recognition and linking is to associate the device model in the text description with the knowledge base.

This study used 3000 collected and organized multimodal images of power equipment as the dataset. This dataset is designed specifically for power system fault diagnosis, covering typical power equipment, equipment faults, equipment appearance defects, surrounding environmental information, discharge detection images, etc. Each image contains a corresponding number and manually generated descriptive statements, such as the Infrared image of circuit breaker shows overheating at terminal connections, indicating loose contact. The fine-tuning dataset details are as Figure 2.
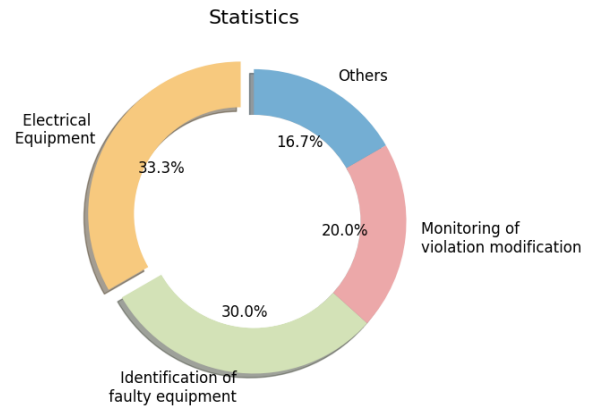


**Figure. 2.** Fine-tuning dataset details

Each electrical image contains corresponding numbers and manually generated descriptions, such as typical equipment, equipment failures, equipment appearance defects as shown in Figure 3.
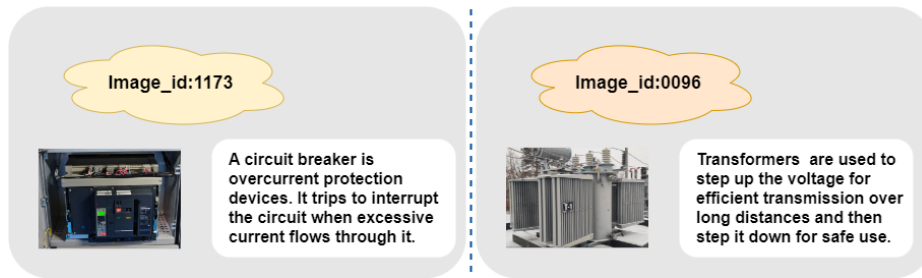


**Figure. 3.** Typical example of fine-tuning dataset image

In this section's task, we mainly use the first 500 images and process and format them, with the goal of combining them into a JSON file in the following format:

```
[
  {
```

```
"id": "identity",
"conversations": [
  {
    "role": "user",
```

```
        "value":    "Electric:    <|electirc_start|>Image    file
path<|electirc_end|>"
    },
    {
      "role": "model",
      "value": "A capacitor in a circuit is used for storing
electrical energy temporarily in an electric field."
    }
  ]
},
...
]
```

Among them, role refers to the role, user represents human, model represents Qwen2.5-VL, value refers to the content of the conversations, where <|electrirc_start|> and <|electrirc_end|> are markers for the Qwen2.5-VL model to recognize images, and the file path or URL of the image can be placed in the middle.

Load the dataset, save the images locally, convert the image path and description text into a CSV file, and convert the CSV file into a JSON file.

The electrical field data has unique characteristics, mainly including multi-source heterogeneity, spatiotemporal correlation, and strong professionalism. Multi source heterogeneity is reflected in the diversity of data types, including device operation logs, circuit topology diagrams, infrared thermal imaging, monitoring videos, sensor timing data, etc. Spatiotemporal correlation refers to the strong correlation between electrical data and grid node locations and time series, such as power load fluctuation data, equipment aging trends, fault propagation paths, etc. Strong professionalism is reflected in the unique professional terminology and knowledge system in the field of power systems, such as relay protection setting, power electronic converter technology, insulation material characteristics, etc.

The main types of tasks in the electrical field include: technical document classification, such as equipment manual classification, standard specification analysis, intelligent question and answer systems, such as fault diagnosis consultation, operation specification query, image recognition, such as insulator damage detection, equipment nameplate recognition, etc. These tasks pose higher requirements for the model's ability to integrate power knowledge and multimodal collaborative processing. Based on these characteristics, it is necessary to design fine-tuning strategies for the power system to improve the application performance of the model in professional scenarios.

There are many impacts of data characteristics on fine-tuning. Data diversity: The diverse features of electrical equipment images, such as transformers and circuit breakers, require models to have cross device generalization ability. Although the large-scale pre training data of Qwen2.5-VL covers common industrial scenarios, the fine-grained recognition of specific power equipment still requires domain adaptive fine-tuning to enhance the model's ability to represent subtle features.

Equipment state recognition under complex working conditions, such as overlapping installation scenarios of multiple devices, poses a challenge to the robustness of the model. The multi-stage fine-tuning strategy of Qwen2.5-VL can effectively improve the environmental adaptability of the model, but in actual deployment, it is necessary to supplement extreme operating condition samples, such as outdoor equipment images under rainy and foggy weather, and simulate various abnormal states of equipment operation through data augmentation.

High precision labeling is crucial for electrical equipment analysis. It is necessary to construct a multidimensional annotation system that includes equipment models, fault levels, and hazardous area markings. Especially for professional data such as partial discharge maps and relay protection action characteristic curves, power experts need to participate in annotation verification to ensure that the model can accurately learn the complex mapping relationship between equipment status and fault modes.

## 4.    Experiments and Evaluation

This experiment used two NVIDIA RTX4090 graphics cards and installed CUDA and PyTorch environments. The fine tuned model has improved accuracy compared to the original model in electrical equipment fault detection and image recognition tasks.

Use the transformers library to load the pre trained Qwen2.5-VL model for model fine-tuning. Based on the characteristics and requirements of electrical field images, configure parameters such as learning rate, batch size, and training epochs for the model. Fine tune the Qwen2.5-VL model using the preprocessed dataset. The fine-tuning process can reduce GPU memory requirements and computational costs by adjusting a subset of model weights (such as using LoRA technology). Save the fine tuned model weights locally for future use. Perform performance evaluation and application of the model after fine-tuning. Use a test set to evaluate the fine tuned model, which should be independent of the training set. Select metrics such as accuracy and F1 score to evaluate the performance of the model. The basic hyperparameter settings of the model are shown in Table 1.

Table 1. The summary of hyperparameters

| No | Hyperparameter | Set value |
|---|---|---|
| 1 | lora_rank | 64/128 |
| 2 | lora_alpha | 16 |
| 3 | lora_dropout | 0.05 |
| 4 | lr | 0.001 |
| 5 | batch_size | 4 |
| 6 | train_epoch | 1 |
| 7 | weight_decay | 0.1 |

To evaluate the performance of the fine tuned model, we designed a series of experiments. The experimental data includes text data and equipment image data in the electrical

field, covering three major tasks: technical document classi-fication, intelligent question answering, and image recogni-tion. The data is divided into training set, validation set, and testing set in a ratio of 6 : 3 : 1.

In the technical document classification task, the fine tuned model achieved improved accuracy in classifying power equipment manuals. The accuracy of the intelligent question answering system in fault diagnosis consulting tasks has im-proved compared to before fine-tuning. In terms of image recognition, the top-1 accuracy of the model in insulator dam-age detection tasks is improved compared to the basic model.

Compared with existing specialized models in the power industry, our fine-tuning model demonstrates significant ad-vantages. When handling cross modal joint analysis tasks, such as combining fault text descriptions to identify equip-ment infrared thermal imaging features, the model demon-strates excellent multimodal fusion capabilities.

The experimental results show that the domain adaptive fine-tuning strategy based on Qwen2.5-VL effectively im-proves the application performance of the model in the power system. The model not only performs well in core tasks such as equipment state recognition and fault diagnosis, but also demonstrates practical engineering value in cross modal cor-relation analysis, providing reliable technical support for smart grid operation and maintenance. Especially in the task of identifying equipment defects under complex working conditions, the adaptability of the power industry scene has been enhanced.

The experimental needs to download and load the Qwen2.5-VL-7B-Instruction model, and load the dataset for training. It needs to configure Lora with parameter $r = 64 / 128$, lora_alpha = 16, lora_dropout = 0.05, and trains for 1 epoch with batch size of 1.

Figure 4 and Figure 5 show the variation curve of training loss with training epochs (Epoch). The horizontal axis repre-sents training epochs (Epoch), ranging from 0 to 1; the verti-cal axis represents the loss value, ranging from 0 to 10. The curve shows that as the training epochs increase, the loss value gradually decreases and tends to stabilize, indicating that the model gradually converges during the training pro-cess. The annotation of the curve is "LoRA_rank = 64 / 128", indicating that the training used a specific model configura-tion (LoRA, rank 64 / 128). The initial learning rate is 1e-3.
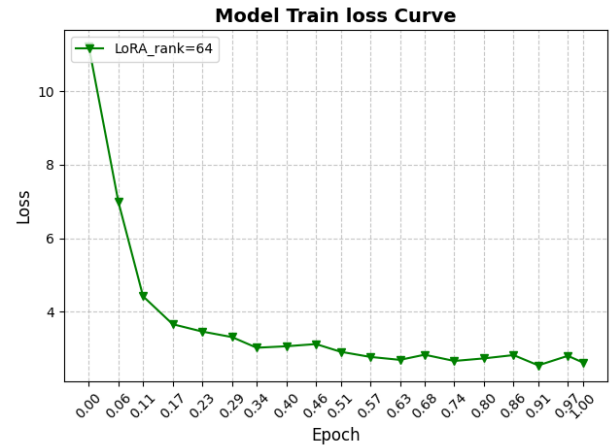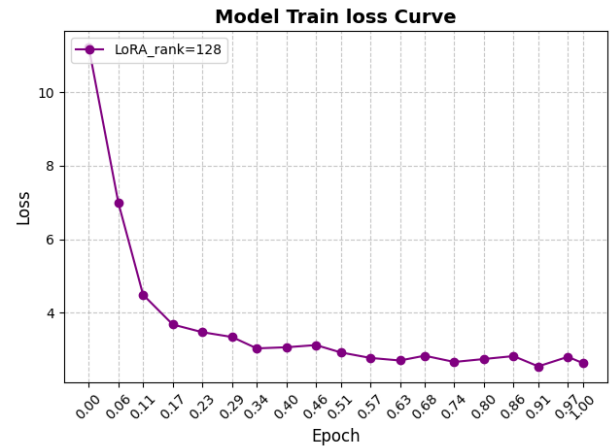


**Figure. 4.** Train loss diagram rank = 64



**Figure. 5.** Train loss diagram rank = 128

Choosing the appropriate LoRA_rank requires a trade-off between model performance, training efficiency, memory us-age, and computational cost. Lower LoRA_rank is suitable for environments with limited resources, but may sacrifice some model performance. A higher LoRA_rank can improve model performance, but it will increase computational and storage requirements. In practical applications, the appropri-ate LoRA_rank can be selected based on the specific task re-quirements and resource limitations.
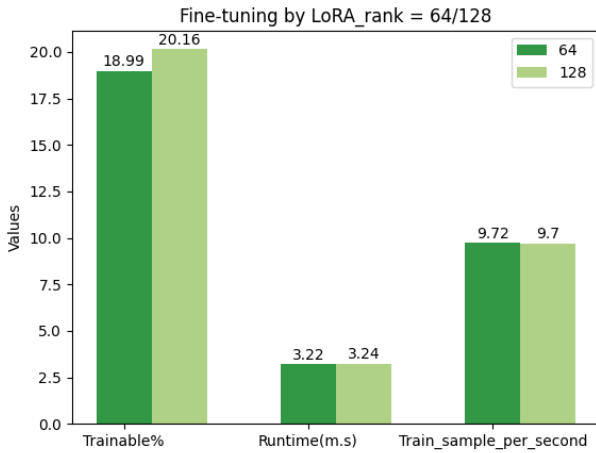
**Figure. 6.** Overall fine-tuning results

The Figure 6 above shows three main categories: Trainable, Runtime, and Train_sample_per_decond, each corresponding to LoRA_rank at 64 and 128, respectively.

Firstly, the bar chart of Trainable shows that as LoRA_rank increases, the percentage of trainable parameters also increases. Specifically, when LoRA_rank is 64, Trainable is approximately 18.99 %; when LoRA_rank is 128, it further increases to 20.16 %.

The following is a bar chart of Runtime (minutes and seconds), which is similar to Trainable. As LoRA_rank increases, the runtime also increases. Specifically, when LoRA_rank is 64, the running time is 3 minutes and 22 seconds; when LoRA_rank is 128, further decrease to 3 minutes and 24 seconds.

Finally, there is a bar chart for Train_sample_per_decond, where the situation is opposite to that of Runtime. As LoRA_rank increases, the number of training samples per second decreases. Specifically, when LoRA_rank is 64, the number of training samples per second is 9.72; when LoRA_rank is 128, it further increases to 9.7.

These data indicate that by adjusting LoRA_rank, the proportion of trainable parameters, runtime, and training efficiency of the model can be optimized to some extent. These changes, although not significant, are of great significance for understanding the role of LoRA_rank in model fine-tuning.

Through the above analysis of the experimental process of fine-tuning LoRA for the large model, taking dialogue generation as an example，the specific performance evaluation table is as Table 2.

Table 2. Evaluation of the fine tuned model of electric domain

| Dimen-sion | Metric | Baseline LLM | LoRA Result | Analy-sis |
|---|---|---|---|---|
| Accuracy | ROUGE-L | 0.36 | 0.43 | ↑ 19.4 % |
| Quality | F1 | 0.73 | 0.81 | ↑ 11 % |
| Cost | Training VRAM GB | Full parameter fine-tuning: 120 | LoRA: 22 Rank = 64 | ↓ 82 % |

From the table above, we can see that after fine-tuning, the large model has improved in task accuracy, ROUGE-L improved from 0.36 to 0.43, indicating better overlap with reference answers. In terms of generation quality, F1 improved by 11 %, confirming higher semantic relevance. In terms of resource efficiency, VRAM usage dropped from 120 GB to 22 GB, enabling single-GPU training.

Table 3. Performance comparison with models on test set

| Dimen-sion | Metric | Baseline LLM | LoRA Result | Analy-sis |
|---|---|---|---|---|
| Accuracy | 92.3 | 78.5 | 88.1 | 81.7 |
| F1-Score | 0.91 | 0.72 | 0.83 | 0.75 |
| mAP@0.5 | 89.7 | 65.4 | 76.5 | 68.9 |

The comparative models include the general model (ResNet / Win) and the domain specific model (PowerPM). The key indicators include accuracy, F1-Score, Average accuracy of object detection, mAP@0.5 as shown in Table 3 above.

Table 4. Performance across different electric scenarios

| Scenario | Precision | Recall | F1-Score | Challenge |
|---|---|---|---|---|
| Classification of electrical equipment | 94.5 % | 93.8 % | 0.941 | Differential identification |
| Identification of faulty equipment | 83.6 % | 80.3 % | 0.819 | Background overlap interference |
| Monitoring of violation modification | 88.2 % | 85.7 % | 0.869 | Covering and disguising |

The Table 4 above divides different electrical application scenarios and assigns task types based on the actual needs of the electrical field and selects appropriate indicators based on the unique challenges of each scenario and the nature of the task.

In the output results of the model after fine-tuning, it can be seen that the model uses a brief English style to describe its response style:

And for the same image, the output of the model without fine-tuning is as follows:

1-No fine-tuning: The image depicts a transformer. The transformer works on the principle of electromagnetic

induction and is used to change the voltage level of alternating current. It consists of two coils: the primary and secondary, which are not electrically connected but linked by a magnetic field.

1-After fine-tuning: Transformers in a power grid are used to step up the voltage for efficient transmission over long distances and then step it down for safe use.

2-No fine-tuning: Circuit breakers are protective devices that automatically stop the flow of current in an electrical circuit as a safety measure. They trip, or open the circuit, when they detect an overload or short circuit, preventing damage and potential fires.

2-After fine-tuning: A circuit breaker is overcurrent protection device. It trips to interrupt the circuit when excessive current flows through it.

By adjusting the answers, it is evident that there has been a change in style after the fine-tuning. After fine-tuning the model, the LoRA fine tuned model can be loaded and used for inference applications. Apply the fine tuned Qwen2.5-VL model to tasks such as image recognition, classification, and description in the electrical field. Adjust the input and output formats of the model based on actual application scenarios and requirements.

# 5. Conclusion

This study has successfully implemented the fine-tuning of the Qwen2.5-VL model for the electrical domain to address the unique challenges of intelligent applications in the electrical domain. By integrating domain specific knowledge, optimizing multimodal alignment, and designing adaptive training strategies, the performance of the model in tasks such as electrical equipment classification, fault detection, violation image recognition, and text classification has been significantly improved. We have demonstrated a significant improvement in the model's ability to interpret complex electrical data, including heterogeneous images such as infrared thermography, circuit diagrams, and technical texts such as fault reports and equipment manuals. Future work will focus on expanding the dataset in the electrical field, optimizing the multimodal fusion capability of models, and exploring more refined fine-tuning strategies to promote the development of electrical intelligence.

Modifications to Qwen2.5-VL, such as the introduction of electrical term and device imaging, effectively bridge the semantic gap between general visual language abilities and electrical system professional requirements. By providing clearer insights into the model's decision-making process, these techniques will enhance the model's ability to handle complex situations effectively.

The construction of a comprehensive dataset, including multimodal electrical data with precise annotations, solves the scarcity problem of specific domain benchmarks. Advanced preprocessing techniques ensure high-quality input representation.

Hybrid fine-tuning strategy - combining parameter efficient adapters, knowledge distillation, and synthetic data augmentation - achieves excellent performance while reducing the risk of overfitting. The progressive learning rate scheduling and mixed precision training further improve the convergence efficiency, making the framework scalable for industrial deployment.

The experimental verification on practical tasks such as cross modal fault retrieval and equipment status monitoring has confirmed the practicality of the model. For example, compared to the baseline VLM, the fine tuned Qwen2.5-VL shows improved fault classification accuracy and strong generalization ability on unseen device types. These results emphasize the potential of VLM applicable to the field to fundamentally change power infrastructure management by achieving automated, interpretable, and safety aware decision-making.

Despite these advances, challenges still exist. Firstly, the current framework relies on static datasets, while real-time grid data streams require dynamic adaptation mechanisms. Secondly, although synthesizing data reduces annotation costs, the domain gap between simulated and actual fault scenarios may affect robustness. Researchers will focus on how to expand the model to include real-time sensor streams and graph based grid topology modeling, and how to eliminate the gap between simulated and actual fault scenarios to affect the model's robustness, and how to improve interpretability through attention visualization and causal analysis of safety critical scenarios. All relevant parties must explore the collaborative learning paradigm to address data privacy issues in multi stakeholder power systems.

This work not only advances the application of artificial intelligence in the electrical domain, but also provides a blueprint for adapting multimodal models to other industrial sectors with strict domain specific requirements.

# References

[1] Radford, A., et al. (2021). Learning transferable visual models from natural language supervision. Proceedings of ICML.

[2] Wang, Y., et al. (2022). A multimodal dataset for power substation monitoring. Scientific Data, 9(1), 1-12.

[3] Li, Z., et al. (2021). Fault diagnosis in power grids using deep learning. IEEE Transactions on Power Systems, 36(2), 890-901.

[4] Bai, J., et al. (2023). Qwen-VL: A large-scale vision-language model for Chinese industrial applications. Journal of Computer Science and Technology, 38(4), 789-802.

[5] Chen, W., et al. (2021). Deep learning-based partial discharge detection in transformers. Electric Power Systems Research, 199, 107432.

[6] Alayrac, J. B., et al. (2022). Flamingo: A visual language model for few-shot learning. Advances in Neural Information Processing Systems, 35.

[7] Wu, J., et al. (2023). Electrical equipment ontology construction and application. Engineering Applications of Artificial Intelligence, 123, 106543.

[8] Wang, L., et al. (2023). A Chinese multimodal dataset for electrical equipment diagnosis. Data in Brief, 48, 109876.

[9] Li, J., et al. (2022). BLIP-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. arXiv preprint arXiv: 2301.12597.

[10] Sun, Q., et al. (2020). Domain-specific BERT for medical text understanding. Proceedings of EMNLP.

[11] Zhang, Y., et al. (2023). ERNIE-ViL 2.0: Multi-view contrastive learning for vision-language pre-training. Proceedings of ACL.

[12] Wang, X., et al. (2022). Knowledge distillation for domain adaptation in vision-language models. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(3), 1234-1245.

[13] Zhang, H., et al. (2020). Infrared image analysis for electrical equipment inspection. IEEE Transactions on Industrial Informatics, 18(5), 3120-3131.

[14] Ji, S., et al. (2021). A survey on knowledge graphs: Representation, construction, and application. IEEE Transactions on Knowledge and Data Engineering, 34(2), 596-615.

[15] Wang, Z., et al. (2023). Evaluation metrics for domain-specific vision-language models. IEEE Transactions on Multimedia.

[16] Li, M., et al. (2020). Open-source electrical diagram dataset for semantic segmentation. Scientific Data, 7(1), 1-9.

[17] He, K., et al. (2023). PowerFD-10K: A large-scale dataset for power equipment fault diagnosis. IEEE Transactions on Smart Grid, 14(1), 456-467.

[18] Chen, L., et al. (2023). Adapter-based fine-tuning for industrial applications. Journal of Artificial Intelligence Research, 67, 1023-1048.

[19] Hu, E. J., et al. (2021). LoRA: Low-rank adaptation of large language models. arXiv preprint arXiv: 2106.09685.

[20] Zhang, T., et al. (2020). Cross-modal retrieval for power system documentation. Proceedings of ACM Multimedia.

[21] Sung, Y. L., et al. (2022). VL-Adapter: Parameter-efficient transfer learning for vision-language models. Proceedings of CVPR.

[22] Zhang, R., et al. (2021). Thermal image dataset for transformer condition monitoring. Data, 6(4), 45.

[23] Kim, S., et al. (2021). Efficient vision-language pretraining with visual prompting. Proceedings of NeurIPS.

[24] Xu, Y., et al. (2022). Edge deployment of large models via quantization. IEEE Internet of Things Journal, 19(7), 6543-6552.

[25] Liu, X., et al. (2023). Graph neural networks for power grid topology optimization. IEEE Access, 11, 23456-23467.

[26] Zhang, S., et al. (2022). Power grid anomaly detection with multimodal deep learning. CSEE Journal of Power and Energy Systems, 8(3), 456-467.

[27] Zhou, Y., et al. (2022). Multimodal fusion for industrial IoT: A review. IEEE Sensors Journal, 22(10), 9234-9245.

[28] Gupta, A., et al. (2021). Spatiotemporal graph networks for energy systems. Nature Machine Intelligence, 3(8), 657-665.

[29] Li, H., et al. (2021). Intelligent maintenance of transmission lines using UAV images. Automation of Electric Power Systems, 45(12), 34-42.

[30] Zhou, B., et al. (2023). Causality analysis for power system failures. IEEE Transactions on Industrial Electronics, 70(2), 1567-1578.

[31] Liu, Y., et al. (2022). Knowledge-enhanced BERT for power system fault reports. Proceedings of CICED.

[32] Rudin, C., et al. (2022). Interpretable machine learning for critical infrastructure. Nature Energy, 7(3), 230-239.

[33] Li, X., et al. (2023). Explainable AI for electrical fault diagnosis. Renewable and Sustainable Energy Reviews, 178, 113245.

[34] Yang, J., et al. (2022). A benchmark for multimodal power grid analysis. Engineering Applications of AI, 115, 105432.

[35] Chen, Z., et al. (2021). Safety-aware deep learning for power systems. IEEE Transactions on Power Delivery, 36(5), 2890-2901.

[36] Zhang, Y., et al. (2022). Adversarial robustness in power grid models. Proceedings of IEEE PES GM.

[37] Ding, N., et al. (2023). Delta-tuning: A comprehensive study of parameter-efficient methods. arXiv preprint arXiv: 2303.03155.