

The Technical Report of DuIVRS-2.

1 EXPERIMENTAL IMPLEMENTATIONS

1.1 Test Data

Table 1 illustrates the detailed statistics of constructed test dataset, which provides a comprehensive evaluation of fine-tuned LLMs from three distinct dimensions.

Table 1. The statistics of test data.

Dataset	$D_{Effectiveness}$	$D_{Generalization}$	$D_{Robustness}$
Number of Samples	626	635	656
Avg. Reply Length	8.52	9.21	21.40
Avg. Reply Frequency	385	193	67

1.2 Training Details

The model training was conducted on Baidu’s PaddleCloud, utilizing eight NVIDIA A100-80G GPUs for fine-tuning. We utilized the AdamW [16] optimizer with the parameters set to $\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 1e - 5$. The training process was configured with a batch size of 128 and a sequence length limited to 1024 tokens. A linear learning rate schedule was applied, incorporating a warm-up phase covering 3% of the total training steps. The maximum learning rates were established at 2×10^{-5} for the EB-turbo model and 1×10^{-4} for the EB-tiny model. For the EB-tiny model, we employed bf16 16-bit (mixed) precision for full-parameter fine-tuning. For the EB-turbo model, we use Lora [9] for parameter-efficient fine-tuning. Finally, we fine-tuned the two models for 2 epochs.

1.3 Power Consumption

Building on prior research [29, 33] and power consumption data for GPU devices, we aim to estimate the financial costs and carbon emissions associated with our training process. Along with previous work, our analysis excludes additional power requirements, such as those from interconnects or ancillary non-GPU energy expenditures. At each iterative stage, the training duration for EB-tiny is about 1 hour, and EB-turbo is 14 hours, amounting to a cumulative $(1 + 14) \times 8 \times 6 = 720$ GPU hours on A100-80G units with a TDP of 400W. Considering GPUs’ actual power use (typically under 400W) and an electricity rate of 1.2 RMB/kWh, the maximum training expense is roughly 400 RMB, with carbon emissions approximating 122kgCO₂eq. Furthermore, utilizing the ERNIE 4.0 API service adds to the cost. With a rate of 0.15 RMB per 1k tokens and across five iterations totaling 72,000 requests at 0.5k tokens each, the API costs about $0.15 \times 72 \times 10^3 = 5400$ RMB. Overall, the complete training costs are projected to stay below 10,000 RMB, with carbon emissions under 1tkgCO₂eq.

2 COOPERATIVE ITERATIVE LEARNING FRAMEWORK

Algorithm 1 delineates the complete cooperative iterative learning framework encompassing data growth and policy improvement phases.

Algorithm 1 Cooperative Iterative Learning framework

- 1: **Input:** \mathcal{D}_g^0 : the original dataset derived by DuIVR-1, $\mathcal{F}(\Theta)$: the parameter of LLM-S, $\mathcal{G}(\Theta)$: the parameter of LLM-L, T: number of iterations.
 - 2: Train $\mathcal{F}(\Theta)$ using Equation 9 on \mathcal{D}_g^0 .
 - 3: **for** iteration t in $\{1, \dots, T\}$ **do**
 - 4: **// Data Growth Phase**
 - 5: Generate dataset \mathcal{D}_g^t according to Equation 10.
 - 6: **for** each sample in \mathcal{D}_g^t **do**
 - 7: Use LLM-L for evaluation according to Equation 15.
 - 8: Use Black-box LLM for evaluation by prompt engineering.
 - 9: Vote and refine according to cooperative voting scheme.
 - 10: **end for**
 - 11: Construct the refined dataset $\tilde{\mathcal{D}}_g^t$, and evaluation dataset $\tilde{\mathcal{D}}_e^t$.
 - 12: **// Policy Improvement Phase**
 - 13: Improve LLM-S by fine-tuning the $\mathcal{F}(\Theta)$ on $\tilde{\mathcal{D}}_g^t$.
 - 14: Improve LLM-L by fine-tuning the $\mathcal{G}(\Theta)$ on $\tilde{\mathcal{D}}_g^t$ and $\tilde{\mathcal{D}}_e^t$.
 - 15: Improve Black-box LLM by optimizing evaluation prompt.
 - 16: **end for**
 - 17: **Output:** Policy $\mathcal{F}(\Theta)$
-

3 PROMPT EXAMPLES

In this section, we showcase various prompts devised for both the training and evaluation phases of DuIVRS-2. Figure 1a provides an example for training the inference model, which adopts a selective generation strategy with CoT mechanism to ensure both safety and stability upon deployment. The model response, “User’s reply confirmed as ABC cake shop” serves as an instance of the chain of thought process, with option “E” indicating the subsequent query “This is Baidu Maps. Are you still operating?”. Figure 1b offers insight into the training evaluation ability of EB-turbo, where it functions to assess the accuracy of EB-tiny’s outputs by yielding a “True/False” verdict.

Input Prompt

You are the Chat Assistant of Baidu Maps, please analyze the conversation, and give the suitable answers. Given that:
POI's name: ABC cake shop
POI's address: {\$address}

The historical conversation:
AI: Hello, is this ABC cake shop?
USER: Yes, who is that?

You need to think and answer step-by-step. First analyze the reply of the user, then given choices of response below, select one that best replies the user.

Choices:
A. This is Baidu Maps confirming information on the Map. Am I reaching ABC cake shop?
B. Sorry, I missed it. Is this ABC cake shop?
C. Have you changed name, or closed?
D. Are you still located at {\$address}?
E. This is Baidu Maps. Are you still operating?
F. Got it. Thank you for supporting Baidu Map, enjoy the rest of the day!
G. Hello, Is this ABC cake shop?
H. Yes, I am an automatic chat assistant, am I reaching to ABC cake shop?
I. Sorry to interrupt, have a good day, bye.

Model Response

User's reply confirmed as ABC cake shop and asked who was calling. E

(a) Example for training the inference model.

Input Prompt

You are the Chat Assistant of Baidu Maps, please analyze the conversation, and give the suitable answers. Given that:
POI's name: ABC cake shop
POI's address: {\$address}

The historical conversation:
AI: Hello, is this ABC cake shop?
USER: Yes, who is that?

The generated response:
Think: User's reply confirmed as ABC cake shop and asked who was calling.
Reply: This is Baidu Maps. Are you still operating?

Does the reply generated above match the historical conversation?

Model Response

True

(b) Example for training the evaluation model.

Fig. 1. Training example for fine-tuning LLMs.

Furthermore, the prompt iteration process for ERNIE 4.0 during the evaluation phase is illustrated in Figure 2. Iteration-V1 enhances the prompt by integrating the criterion “expresses a wish to hang up”, and iteration-v2 further refines it by

adding “indirectly expresses a wish for the conversation to end”, thus extending coverage to more domain-specific edge cases and enhancing the discernment precision of ERNIE 4.0.

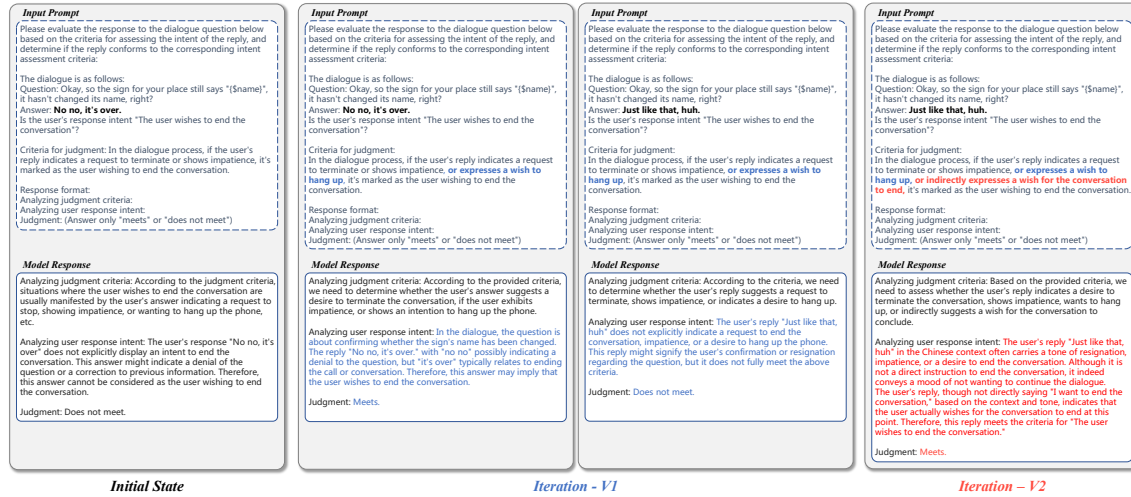


Fig. 2. The prompt iteration of ERNIE 4.0 in evaluation stage.