# OpsEval: A Comprehensive Benchmark Suite for Evaluating Large Language Models' Capability in IT Operations Domain

**Yuhe Liu**[1]   **Changhua Pei**[2]   **Longlong Xu**[1]   **Bohan Chen**[1]
**Mingze Sun**[1]   **Zhirui Zhang**[3]   **Yongqian Sun**[4]   **Shenglin Zhang**[4]
**Kun Wang**[1]   **Haiming Zhang**[2]   **Jianhui Li**[2]   **Gaogang Xie**[2]
**Xidao Wen**[5]   **Xiaohui Nie**[2]   **Minghua Ma**[6]   **Dan Pei**[1]
[1]Tsinghua University   [2]Chinese Academy of Sciences
[3]Beijing University of Posts and Telecommunications   [4]Nankai University
[5]BizSeer   [6]Microsoft

```
lyh23@mails.tsinghua.edu.cn, chpei@cnic.cn,
{xull23,cbh22,smz19}@mails.tsinghua.edu.cn, wannabtl@bupt.edu.cn,
{sunyongqian,zhangsl}@nankai.edu.cn, wangkun_cs@tsinghua.edu.cn,
{hai,lijh,xie}@cnic.cn, wenxidao@tsinghua.edu.cn, xhnie@cnic.cn,
minghuama@microsoft.com, peidan@tsinghua.edu.cn
```

## Abstract

The past decades have witnessed the rapid development of Information Technology (IT) systems, such as cloud computing, 5G networks, and financial information systems. Ensuring the stability of these IT systems has become an important issue. Large language models (LLMs) that have exhibited remarkable capabilities in NLP-related tasks are showing great potential in AIOps, such as root cause analysis of failures, generation of operations and maintenance scripts, and summarizing of alert information. Unlike knowledge in general corpora, knowledge of Ops varies with the different IT systems, encompassing various private subdomain knowledge, sensitive to prompt engineering due to various sub-domains, and containing numerous terminologies. As a consequence, existing NLP-related benchmarks such as C-Eval and MMLU can not guide the selection of suitable LLMs for Ops, and current metrics like BLEU and ROUGE can not adequately reflect the question-answering (QA) effectiveness in the Ops domain. Therefore, this paper proposes a comprehensive benchmark suite named **OpsEval**, including an Ops-oriented evaluation dataset, an evaluation benchmark for Ops, and a specially designed QA evaluation method for Ops. Our dataset contains 7,184 multiple-choice questions and 1,736 QA questions. We have carefully selected and released 20% of the dataset by CC-BY-NC-4.0 license[1] written by domain experts in various sub-domains to assist current researchers in preliminary evaluations of their LLMs tailored for Ops (OpsLLM). The remaining undisclosed 80% of the data is used to prevent test set leakage. We test over 21 latest LLMs under various settings such as self-consistency, chain-of-thought, and in-context learning, revealing findings when applying LLMs to Ops. We also propose an evaluation method for QA in Ops, which has a coefficient of 0.9175 with human experts and is improved by 0.2470 and 1.313 compared to BLEU and ROUGE, respectively. Over the past seven months, our dataset and leaderboard[2] have been continuously updated. For reproducibility, the evaluation framework code is publicly available.

---

[1]Data page is available at https://github.com/NetManAIOps/OpsEval-Datasets.
[2]Leaderboard is available at https://opseval.cstcloud.cn/content/leaderboard.

# 1 Introduction

The IT Operations (Ops) plays a crucial role in maintaining the efficient and stable operation of information systems such as cloud computing, 5G networks[3] and financial information systems. As the Internet continues to expand rapidly, the scale and complexity of systems are escalating, leading to the emergence of artificial intelligence-assisted operations as a novel trend. Termed "AIOps" by Gartner (Lerner, 2017), this technique utilizes artificial intelligence to address (but is not limited to) tasks such as anomaly detection, fault analysis, generation of alert summaries, performance optimization, and capacity planning.

In recent years, large language models (LLMs) have witnessed significant advancements. The latest models, such as GPT-4o (OpenAI, 2024), GPT-4V (OpenAI, 2023b), Meta-Llama-3 (AI@Meta, 2024), and GLM-4 (Zeng et al., 2022), have demonstrated exceptional generalization and task-planning capabilities. As a result, these models have provided numerous opportunities to enhance downstream domain-specific applications. With its advanced summarizing, report analyzing, and ability to diagnose errors, LLM is well suited for Ops on tasks like question answering, information summarizing, and report analysis. Hereinafter, we refer to the LLM used for Ops as **OpsLLM**, regardless of whether they have been optimized specifically for Ops.

While there are benchmarks for assessing general-purpose NLP-related capabilities, such as C-EVAL (Huang et al., 2023) and AGIEval (Zhong et al., 2023), there are also benchmarks for specific domains, like FinEval (Zhang et al., 2023) in the financial sector and CMB (Wang et al., 2023a) in the medical sector. However, no benchmark exists to evaluate the effectiveness of LLMs or OpsLLMs in Ops tasks. There is an urgent need for an Ops benchmark that informs us about the performance of current LLMs on Ops tasks. On the other hand, a good benchmark can significantly aid the optimization process of OpsLLMs tailored for the Ops domain. Nevertheless, due to the specialty of the Ops tasks, constructing an Ops benchmark presents the following challenges:

- **Sensitive data.** The Ops data is primarily sensitive and proprietary to companies, with very few publicly available data, making it difficult for any company to independently provide sufficient evaluation data to ensure confidence in the test results.

- **Sub-domains.** There are many different sub-domains, such as 5G communications, cloud computing, and bank transaction systems. Ops within each sub-domain typically requires a combination of various capabilities, such as network configuration or terminology explanation, which we refer to as "tasks". Due to the large number of sub-domains and tasks in the Ops field and the lack of an authoritative and systematic taxonomy, classifying a large number of questions becomes a challenge.

- **Prompt sensitivity.** Due to the relatively proprietary nature of the Ops, existing LLMs have not undergone specialized supervised fine-tuning (SFT) for instruct following within the Ops field, the evaluation results are more sensitive to prompt engineering. Designing appropriate prompts for robust and accurate evaluation is challenging.

- **QA metric.** Existing metrics like BLEU only consider the similarity of model output to the reference in natural language aspects, which does not always reflect true performance in Ops tasks. Some terms and expressions in Ops sub-domains have specific meanings that LLM cannot summarize. Designing an automatic evaluation method that assesses the accuracy of QA of Ops from an accurate semantic level is challenging.

To address these issues, we propose **OpsEval**, a comprehensive benchmark suite for evaluating LLMs' capability in the IT operations domain. First, to tackle the challenge of benchmark data mostly being private and not publicly shareable, we initiated a community around AIOps, which has attracted dozens of companies to participate. We have selected 10 representative sub-domains from the community, allowing continuous data contributions from community members. We then aggregate data under the same sub-domain to ensure robustness in evaluation. Additionally, we generated multi-choice (MC) and question-answering (QA) questions as supplements based on publicly available network management books. To address the challenge of classifying the numerous sub-domains and tasks in the Ops field, we employ model-based pre-clustering and manual review to

---

[3]Strictly speaking, 5G belongs to the field of communications technology (CT), but given its broad association with the information technology (IT) sector, for the sake of generality, we refer to it as IT operations, abbreviated as Ops, throughout the remainder of this paper.

Table 1: A comparison of OpsEval with other popular datasets/benchmarks.

| | MMLU | HELM | BIG-bench | SEAL | C-Eval | AGIEval | FLUE | MultiMedQA | CMB | NetOps | OpsEval |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Ops Domain Dataset | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Open-sourced Benchmark | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| Up-to-date Leaderboard | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |

annotate eight tasks and three abilities for independent evaluation. Considering the prompt sensitivity of benchmark results, we systematically test model performance under self-consistency (SC), chain-of-thought (CoT), and few-shot in-context learning (ICL). The prompts used in our evaluation are also disclosed in the paper. Lastly, to address the inaccuracy of existing metrics in Ops QA evaluation, we design FAE-Score, which assesses questions from the perspectives of fluency, accuracy, and evidence. Experimental results show that our designed metrics align strongly with the annotations of human experts.

The contributions of our paper are as follows: **1)** We introduce **OpsEval**, the first bilingual multi-task dataset in the operations domain, covering 8 tasks and 3 abilities with 8920 questions. To assist researchers in preliminary evaluating their OpsLLMs, we have carefully selected and released 20% of QAs from our benchmark licensed under CC-BY-NC-4.0, with the remaining 80% of undisclosed data preventing unfair evaluations due to data leakage (Wei, et.al., 2023) **2)** Based on the dataset, we introduce the OpsEval evaluation benchmark, conducting independent and robust evaluations with various prompting techniques and a specifically designed evaluation metric named FAE-Score. Compared to the commonly employed BLEU and ROUGE metrics, FAE-Score exhibits a more pronounced congruence with the evaluations of human experts. Specifically, FAE-Score attains a correlation coefficient 0.9175 with expert assessments, surpassing the coefficients of 0.6705 for BLEU and -0.3957 for ROUGE. **3)** We released an online leaderboard that continuously updates mainstream LLMs' performance on Ops tasks for the past 7 months. To ensure the reproducibility and reliability of this leaderboard, we made our evaluation framework's code publicly available[4]. **4)** Based on the results of OpsEval evaluation, we provide key observations and practical lessons to help domain practitioners make decisions such as whether existing models are sufficiently applicable within a specific sub-domain, the necessity for fine-tuning and whether model quantization compromises the effectiveness.

## 2 Related Works

As LLMs evolve rapidly, their complex and varied capabilities are increasingly recognized. As a result, there is a growing trend towards evaluation benchmarks tailored specifically for LLMs. These can be divided into two categories: general ability benchmarks and domain-specific benchmarks.

**General ability benchmarks** assess the general abilities of LLMs across various tasks. These tasks evaluate LLMs' capacity for logical reasoning, general knowledge, common sense, and other similar abilities rather than being confined to a particular domain. MMLU (Hendrycks et al., 2021) is a benchmark designed to measure knowledge acquired during pretraining by evaluating models exclusively in zero-shot and few-shot settings, covering 57 subjects across STEM. HELM (Liang et al., 2022) employs seven distinct metrics in 42 unique scenarios, offering a comprehensive evaluation of LLMs' capabilities across multiple dimensions. BIG-bench (Srivastava et al., 2022) comprises 204 tasks spanning a wide array of topics, with a particular focus on tasks deemed beyond the reach of current LLMs. SEAL (AI, 2024b) features private, expert evaluations of leading frontiers models. C-Eval (Huang et al., 2023) is a comprehensive Chinese evaluation suite designed to assess Chinese LLMs' advanced knowledge and reasoning abilities rigorously. AGIEval (Zhong et al., 2023) curates authentic questions from examinations such as the Chinese College Entrance Exam (CCEE) and the SAT, constructing a fundamentally human-centric evaluation dataset.

**Domain-specific benchmarks** evaluate the abilities of LLMs to handle tasks in specific fields. These benchmarks require LLMs to possess specialized knowledge in a specific domain and to respond in a manner consistent with the cognitive patterns of that field. Despite the rapid progression of LLMs in specialized domains, the evaluation metrics for these specific areas have received less attention. FLUE (Shah et al., 2022) is an open-source comprehensive suite of benchmarks, including new benchmarks across 5 NLP tasks in financial domain. MultiMedQA (Singhal et al., 2022) is

---

[4]Evaluation framework code is available at https://github.com/NetManAIOps/OpenCompass-OpsQA.

Table 2: The number of questions in OpsEval, grouped by their sub-domains.

| Sub-domain | Source | License | Type | Questions |
|---|---|---|---|---|
| Wired Network | Operation Textbooks | Public | Multi-Choice | 3901 |
| 5G Communication | Certification Exams | Public | Multi-Choice | 2615 |
| | | | Question-Answering | 1162 |
| Oracle Database | Company Materials | CC-BY-NC-4.0 | Multi-Choice | 497 |
| Log Analysis | Company Materials | CC-BY-NC-4.0 | Question-Answering | 420 |
| DevOps | Automated Generation | CC-BY-NC-4.0 | Question-Answering | 154 |
| Securities Information System | Company Materials | CC-BY-NC-4.0 | Multi-Choice | 91 |
| Hybrid Cloud | Company Materials | CC-BY-NC-4.0 | Multi-Choice | 40 |
| Financial IT | Company Materials | CC-BY-NC-4.0 | Multi-Choice | 40 |

an extensive medical question-answering dataset, with questions derived from professional medical exams, research, and consultation records. CMB (Wang et al., 2023a) includes multi-choice questions (CMB-Exam) and complex clinical questions based on real case studies (CMB-Clin). NetOps (Miao et al., 2023) focuses on evaluations in the network field, which is relevant to the field of Ops. NetOps includes multi-choice questions in both English and Chinese and a few cloze tests and question-answering questions. However, they only focus on wired network operations and while the dataset is released, they lack a benchmark that continuously updates the leaderboard.
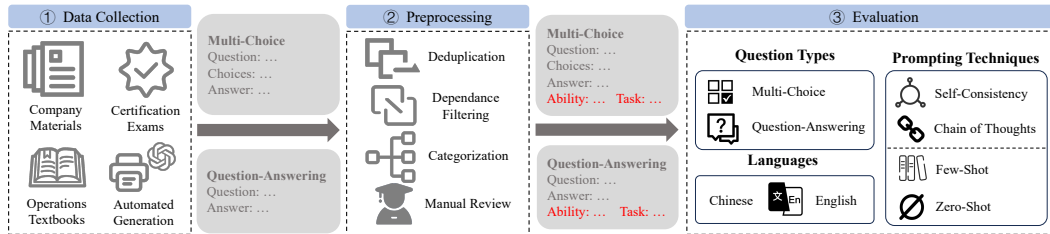
# 3 OpsEval Benchmark



Figure 1: The framework of OpsEval

Figure 1 shows the overall framework of OpsEval from construction to evaluation. We collected data from multiple sources and then preprocessed it to enhance its quality. Finally, we evaluated LLMs on the dataset using various prompt engineering techniques.

## 3.1 Data Collection

Our benchmark questions have been collected from various sources; we summarize them into four categories: company materials, certification exams, Ops textbooks, and automated generation. Each source is highly esteemed globally and reviewed by our Ops collaborators.

**Company Materials.** include production environment materials like Ops tickets and error logs , as well as internal documents and tests for Ops staff training. We have established cooperative relationships with 10 companies, covering various sectors like telecommunications, finance, and Ops service/tool providers, and received expert collaboration and Ops materials from them. The Appendix A.1 provides information about the companies and experts.

**Certification Exams.** include knowledge assessments necessary for becoming an Ops staff and are naturally in the form of multiple-choice and question-answering questions. We obtained the relevant study guidebooks for these certification exams from public book websites and extracted sample questions from them as one of the sources for Ops questions.

**Operations Textbooks**. We first constructed a seeding keyword list for the Ops field and searched for related books. The textbooks contain relatively complete knowledge content, which can provide experts with materials for question creation, and some books themselves also include a certain number of exercises at the end of the chapters.

**Automated Generation.** To enhance the diversity and depth of our test set, we source QAs from authoritative books covering a range of Ops domains by extracting textbook contents and asking GPT4 to generate questions. In Appendix A.7, we discussed the methodology and challenges.

## 3.2 Preprocessing

We systematically carried out the preprocessing of our original data in the following stages:

**Deduplication:** Any repeated or highly similar questions are identified and removed to avoid redundancy in the test set. We calculate the cosine similarity of the question stems by bge-large-zh-v1.5 (Xiao et al., 2023) to detect duplicate questions and identify pairs of questions with a similarity above a certain threshold (th=0.7).

**Dependance Filtering:** We have filtered out questions that rely on external images or document content to ensure the completeness of the question content itself. The filtering process was done by two parallel lists of empirical keywords in the question stems and the responses of GPT-3.5-turbo. The keyword list can be found in the Appendix A.2.

**Question Categorization:** We devise a categorization that captures many tasks that professionals confront in practical applications. The categorization process consists of two steps: automated screening and manual review. We first use GPT-4 for topic modeling to gain rough insights about the dataset and determine the relevance of each question to Ops, which resulted in more than 20 tasks but had an imbalanced distribution[5]. We then involved dozens of experts during the manual review process to categorize the questions into eight tasks and three abilities. The distribution of the questions across these eight tasks and three ability levels is shown in Table 3, and the details of each task and ability can be found in Appendix A.3.

**Manual Review:** In the manual review step, we asked Ops experts from the industry to inspect the results of the previous three automated steps, including confirming duplicate and invalid questions and examining the classification results of GPT-4. In our work, an expert is defined as an individual with ten or more years of professional experience in their field, whether as an employee or a researcher. Experts were also asked to drop the questions unrelated to the Ops field. We split the dataset by n-folds and ensure each fold has at least two experts to review. As listed in Table 2, this quality enhancement process resulted in a refined test set of approximately 7,000 multi-choice and 2,000 question-answering questions.

Table 3: The distribution of different tasks and abilities of questions in OpsEval.

|  | Category | Percentage (%) |
|---|---|---|
| **Task** | Automation Scripts | 3.3 |
|  | Monitoring and Alerting | 5.2 |
|  | Performance Optimization | 5.3 |
|  | Software Deployment | 7.9 |
|  | Fault Analysis and Diagnostics | 13.7 |
|  | Network Configuration | 29.0 |
|  | General Ops Knowledge | 20.2 |
|  | Miscellaneous | 15.5 |
| **Ability** | Knowledge Recall | 49.8 |
|  | Analytical Thinking | 39.9 |
|  | Practical Application | 10.2 |

## 3.3 Evaluation Settings

**Multi-choice questions** offer a structured approach with definitive answers. These questions are straightforward and provide a clear metric for assessment. We use **accuracy** as the metric. A choice-extracting function based on regular expressions is used to extract the predicted answer of LLMs. Then, we calculate the accuracy based on the extracted answer and the ground-truth labels.

**Question-answering questions** do not come with predefined options. We use two metrics for question-answering questions: one is based on word overlaps, and the other is based on semantic similarity. For the first type, we use ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002). For the second type, we use LLM judges and human experts to evaluate the outputs of LLMs, called FAE-Score and Expert-Evaluation, designed explicitly in OpsEval. For these two metrics, we design three criteria highly related to Ops' needs. The three criteria in consideration are as follows:

- **Fluency**. Assessment of the linguistic fluency in the model's output and compliance with the question-answering question's answering requirements.

- **Accuracy**. Evaluation of the precision and correctness of the model's output, including whether it adequately covers key points of the ground-truth answer.

---

[5]For prompts used during the preprocessing, please see in Appendix.

- **Evidence**. Examine whether the model's output contains sufficient argumentation and evidential support to ensure the credibility and reliability of the answer.

For FAE-Score, we provide the judge model with the question and the reference answer, followed by one anonymous model's output. We use GPT-4 as the judge model of FAE-Score based on evidence that GPT-4 can reliably match human preferences with over 80% agreement (Zheng et al., 2023). In the subsequent validation section §5, we also verified the consistency between GPT-4's scores and expert scores. Figure 2b shows the prompt we used to ask the judge model to score. For Expert-Evaluation, we asked experts to score it between 0 and 3 for each criterion. During the scoring, the raw question, the detailed answer and its key points, and the output of an anonymous model are given at each iteration. For the demographical information of the experts, see Appendix A.1.

> I need your help in analyzing a multi-choice question, determine the domain and the task type it belongs to.
> **Domains:** When classifying the domain, be specific, dive deeper into domains such as: Database/Network Operations
> **Task Types:** For the task type, consider categories like: Monitoring and Alerts, Performance Optimization
> **Summary your response as JSON format: {"domain": "specific_domain", "task": "specific_task_type"}**

(a) The prompt for GPT-4 initial categorization

> **{question}**
> Answer: **{keypoint}**
> Explanation: **{reference}**
> Output of a model: **{prediction}**
> Give a score between 1 to 10 according to the answer and explanation, 1 being the least effective and 10 being the most effective, please consider the following three aspects when rating:
>
> - **Fluency**. Assessment of the linguistic fluency in the model's output and compliance with the question-answering question's answering requirements.
> - **Accuracy**. Evaluation of the precision and correctness of the model's output, including whether it adequately covers key points of the ground-truth answer.
> - **Evidence**. Examine whether the model's output contains sufficient argumentation and evidential support to ensure the credibility and reliability of the answer.

(b) The prompt for GPT-4 Score

Figure 2: Prompts for ChatGPT used in the framework of OpsEval.

**Prompting Techniques.** We use various settings to evaluate LLMs on OpsEval to get a comprehensive overview of their performance. We evaluate LLMs in zero and few-shot (3-shot) settings. For each setting, we evaluate LLMs in four sub-settings of prompt engineering, that is, naive answers (Naive), self-consistency (SC) (Wang et al., 2023b), chain-of-thought (CoT) (Wei et al., 2023), self-consistency with chain-of-thought (CoT+SC). We set the number of queries in SC to **5**.

**Models.** We evaluate popular LLMs covering different weights from different organizations. The model selection was guided by specific criteria: We aimed to include the latest and most advanced large language models, with a particular focus on those capable of handling Chinese input. The detailed information of all 21 LLMs can be found in Table 9 in Appendix B.1.

**Fairness Consideration** We are currently making 20% of the data available to the public for Ops community contribution and research purposes, yet for fairness of the evaluation, the complete version of the OpsEval dataset is kept undisclosed. To evaluate a new model, users can submit a Docker image with an initialization script when starting a container based on it. We will run the evaluation automatically and obtain the result on the OpsEval website. Users can choose to disclose their results on the leaderboard of OpsEval or not based on their preference.

This released 20% dataset serves as sample questions so that the researchers will know the types and topics of questions and answers expected by the benchmarks. This allow researchers to gain intuition and insights on how to improve their models. Model developers can use the 20% dataset for locally evaluating by themselves their model's performance, enabling faster iterations of their model training. The 20% dataset can be used as a seed to generate QA pairs using automatic QA generation algorithms (Wang et al., 2023c), ultimately providing more Ops text data for improving new models.
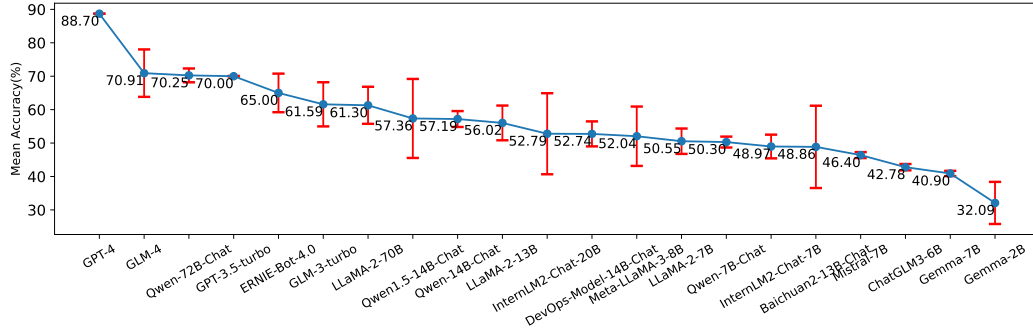
Figure 3: LLMs' overall performance on Wired Network Operations English test set (3-shot). Models are ranked based on their mean accuracy among different settings. The error bars represent the variance in the model's accuracy across different prompting techniques.

Table 4: LLMs' overall performance (Accuracy%) on Wired Network Operations English test set (3-shot). Models are ranked based on their best performance among different settings.

| Model | Naive | SC | CoT | CoT+SC | Best Score | Variance |
|---|---|---|---|---|---|---|
| GPT-4 | / | / | 88.70 | / | 88.70 | / |
| Claude-3-Opus | 79.70 | 78.36 | 81.24 | 82.28 | 82.28 | 2.9638 |
| Qwen2-72B-Instruct | 76.91 | 76.91 | 80.12 | 80.65 | 80.65 | 4.0720 |
| Qwen2-7B-Instruct | 65.74 | 65.74 | 67.42 | 67.42 | 67.42 | 0.9408 |
| GLM-4 | 64.77 | 64.77 | 77.06 | 77.06 | 77.06 | 50.3480 |
| GPT-3.5-turbo | 68.30 | 68.30 | 70.90 | 72.50 | 72.50 | 4.2800 |
| Qwen-72B-Chat | 70.32 | 70.32 | 70.13 | 70.22 | 70.32 | 0.0083 |
| ERNIE-Bot-4.0 | 60.00 | 60.00 | 70.00 | 70.00 | 70.00 | 33.3333 |
| LLaMA-2-70B | 55.00 | 56.20 | 66.80 | 67.20 | 67.20 | 43.5866 |
| Qwen1.5-14B-Chat | 52.23 | 53.52 | 59.53 | 64.17 | 64.17 | 23.0368 |
| DevOps-Model-14B-Chat | 63.85 | 61.96 | 41.15 | 44.01 | 63.85 | 139.6604 |
| GLM-3-turbo | 59.53 | 59.53 | 63.65 | 63.65 | 63.65 | 5.6581 |
| Qwen-14B-Chat | 62.60 | 59.70 | 50.58 | 55.88 | 62.60 | 26.9921 |
| Meta-LLaMA-3-8B | 41.03 | 42.07 | 62.45 | 62.62 | 62.62 | 110.2314 |
| LLaMA-2-13B | 53.30 | 53.00 | 56.80 | 61.00 | 61.00 | 13.9758 |
| InternLM2-Chat-20B | 60.48 | 60.48 | 45.10 | 45.10 | 60.48 | 78.8481 |
| LLaMA-2-7B | 48.20 | 46.80 | 52.00 | 55.20 | 55.20 | 14.4366 |
| Qwen-7B-Chat | 52.10 | 51.00 | 48.30 | 49.80 | 52.10 | 2.6600 |
| Baichuan2-13B-Chat | 51.90 | 51.60 | 44.50 | 47.45 | 51.90 | 12.5822 |
| Gemma-7B | 30.24 | 30.24 | 51.56 | 51.56 | 51.56 | 151.5141 |
| InternLM2-Chat-7B | 48.2 | 48.2 | 49.74 | 49.74 | 49.74 | 0.7905 |
| Mistral-7B | 47.22 | 47.22 | 45.58 | 45.58 | 47.22 | 0.8965 |
| ChatGLM3-6B | 42.10 | 42.10 | 43.47 | 43.47 | 43.47 | 0.6256 |
| Gemma-2B | 26.63 | 26.63 | 37.54 | 37.54 | 37.54 | 39.6760 |

**API Exposure.** For large language models, user requests constitute private data that should be carefully protected. We believe that ensuring data privacy and not using such data for subsequent training is a fundamental requirement that all publicly accessible models must meet. The APIs we utilized for model evaluation explicitly guarantee that user data will not be used for model training (Microsoft, 2023; BigModel, 2023; Cloud, 2023). When evaluating closed-source models such as ChatGPT and GLM, the questions in this non-disclosed questions were sent (thus exposed) to their APIs, but the answers to these questions were still non-disclosed.
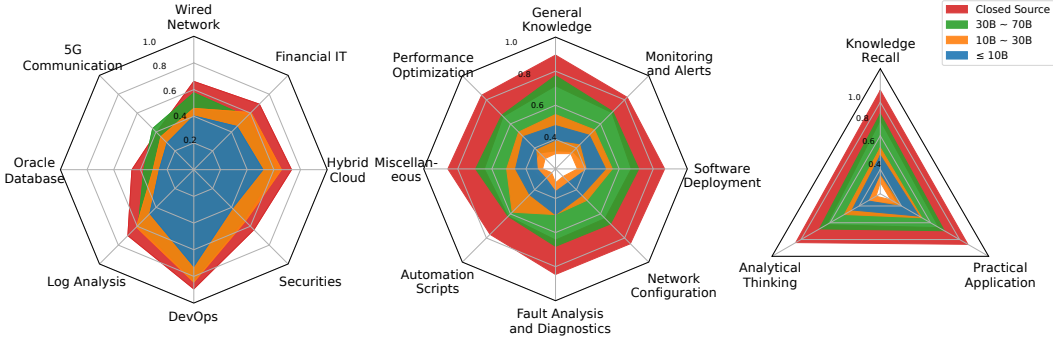
Figure 4: LLMs' performance on eight Ops sub-domains, eight tasks and three abilities. Each colored area presents the lower and upper bound of the corresponding parameter-size group.

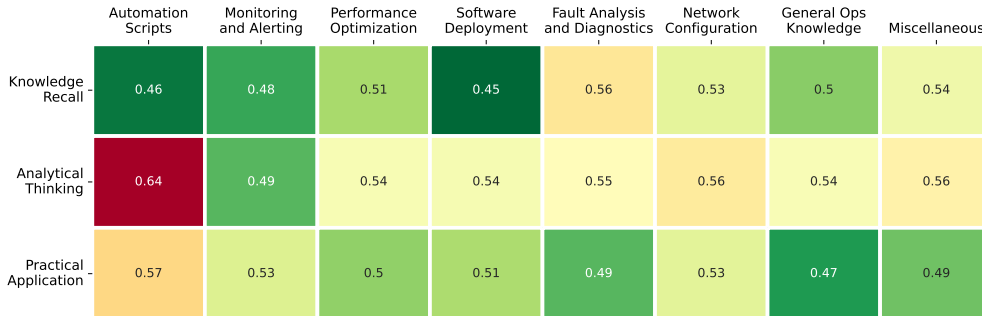| | Automation Scripts | Monitoring and Alerting | Performance Optimization | Software Deployment | Fault Analysis and Diagnostics | Network Configuration | General Ops Knowledge | Miscellaneous |
|---|---|---|---|---|---|---|---|---|
| Knowledge Recall | 0.46 | 0.48 | 0.51 | 0.45 | 0.56 | 0.53 | 0.5 | 0.54 |
| Analytical Thinking | 0.64 | 0.49 | 0.54 | 0.54 | 0.55 | 0.56 | 0.54 | 0.56 |
| Practical Application | 0.57 | 0.53 | 0.5 | 0.51 | 0.49 | 0.53 | 0.47 | 0.49 |

Figure 5: Heatmap of failure case distribution regarding tasks and abilities. The values represent the proportion of failure cases across all LLMs; redder areas indicate higher failure rates.

# 4 Result Analysis

## 4.1 Overall Performance

The results of the few-shot evaluation with four settings on the Wired Network Operation test set are shown in Figure 3. Results of the other sub-domains and settings are shown in Appendix B.4. [6] Overall, open-sourced LLMs yield evaluation results on the OpsEval benchmark generally worse than those in general domains like MMLU (Hendrycks et al., 2021) and CEval (Huang et al., 2023). This comparison highlights the necessity of explicitly fine-tuning OpsLLM for the Ops field. Closed source models like GPT-4 and GLM-4 consistently outperform open source models, while smaller models, such as Qwen1.5-14B-Chat, exhibit competitive performance in multi-choice questions, thanks to their fine-tuning process and the quality of their training data. However, their large variance across the four settings suggests that such models may have worse robustness under different prompts. Furthermore, we observed significant variability in how different LLMs respond to various prompt engineering techniques. Given the critical importance of stability in the Ops domain, it is essential to consider a model's sensitivity to prompts when selecting foundation model. Further research into prompt engineering is needed to improve model performance and reliability in this domain.

**Observations:** 1) Few-shot and CoT can significantly increase performance if the model is tuned to adapt to these techniques, while SC may have little influence on highly consistent LLMs. 2) Smaller models with weaker natural language abilities are less stable with advanced prompts. Simpler prompts work better for them.

**Pratical Lesson**: The choice of fundamental models should be a balance between their performance (average score) and robustness (variance) under different prompt settings.

## 4.2 Performance on Different Tasks and Abilities

To investigate how LLMs perform in each Ops sub-domain and each task, and to what extent they possess the general abilities, we summarize the result of different parameter-size groups of LLM

---

[6]Due to the consideration of time, cost, and API rate limits, for GPT-4, we only make the 3-shot evaluation with the CoT setting to serve as an upper bound of all LLMs to provide a reference.

Table 5: LLMs' performance on English network operations question-answering problems. Total is the sum of the previous three columns.

| Model | ROUGE | BLEU | FAE-Score | Expert Evaluation | | | |
|---|---|---|---|---|---|---|---|
| | | | | Fluency | Accuracy | Evidence | Total |
| GPT-3.5-turbo | 12.26 | 6.78 | 8.47 | 3.00 | 1.96 | 1.20 | 6.16 |
| LLaMA2-70B | 7.74 | 4.20 | 7.28 | 2.92 | 1.48 | 1.32 | 5.72 |
| LLaMA2-13B-Chat | 4.98 | 3.43 | 7.16 | 2.82 | 1.34 | 1.62 | 5.78 |
| Baichuan2-13B-Chat | 4.76 | 0.35 | 5.85 | 2.40 | 1.12 | 1.02 | 4.54 |
| Qwen-7B-Chat | 11.82 | 4.33 | 5.63 | 2.56 | 1.14 | 0.84 | 4.54 |
| ChatGLM3-6B | 9.71 | 5.07 | 4.88 | 2.84 | 0.76 | 0.76 | 4.36 |
| InternLM2-7B-Chat | 13.27 | 0.54 | 4.52 | 1.80 | 0.70 | 0.10 | 2.60 |

and plot them on three radar charts in Figure 4. Regarding the eight tasks we tested, LLMs yield higher accuracy in General Knowledge tasks, while their performance drops and varies drastically in highly specialized tasks like Automation Scripts and Network Configuration, reflecting the impact of specialized corpus and domain knowledge on the performance of LLMs. By grouping LLMs by their parameter size, we find that although LLMs with 10B-30B parameters have higher accuracy in their best cases compared with LLMs with no more than 10B parameters, different 10B-20B LLMs' performance varies drastically. To provide systematic practical lessons for researchers in the operations domain on pre-training and fine-tuning OpsLLM, we have analyzed the error rates of LLMs across the 8 tasks and 3 abilities in Figure 5. By examining the focus areas across different categories, we have identified key research targets for capability training.

**Observations:** Among the 24 categories of results, models performed the worst in Analytical Thinking for Automation Scripts. This indicates that current models can only recall the learned scripts but struggle to infer their logical relationships. Similarly, Analytical Thinking showed the lowest performance across the three major tasks, indicating that current OpsLLM models still have some way to go before becoming foundational models for Ops Agents. Thus, researchers should focus on inference-related SFT (supervised fine-tuning) datasets.

**Insights:** 1) Among different sub-domains of Ops, 5G communication and database demand further pretraining and fine-tuning of LLM. 2) To be capable of an Ops agent, the foundation model must be able to make a connection between specialized domain knowledge.

### 4.3 Performance on Question-Answering

Table 5 presents the evaluation results of 200 question-answering English questions across four metrics: ROUGE, BLEU, FAE-Score, and Expert-Evaluation, sorted by FAE-Score results. To gain more insight into how different metrics perform in QA evaluation, we use Figure 16 (see in Appendix B.8.2) as a case analysis. While BLEU and ROUGE are efficient in natural language comparison, they lack semantic information to determine which part of the context is more important than others. Knowing that a given benchmark evaluates QA based on BLEU/ROUGE, there is an obvious way to trick the metric: repeat patterns occurring in the question, gaining a higher possibility to match some patterns in the reference answer. Due to their lack of semantic information related to Ops and the potential hack, traditional metrics like BLEU are unsuitable for specialized benchmarks. Instead, with specialized prompting, FAE-Score can pick up those important keywords and not be influenced by repeated words that contain no useful information. In a later section, we discuss the alignment between different metrics and expert evaluation, validating the effectiveness of FAE-Score in automated QA evaluation within the Ops domain.

**Practical Lesson:** FAE-Score is suitable for large-scale qualitative evaluations in the Ops field.

### 4.4 Performance on Different Quantization parameters

We conducted experiments on different quantized versions of LLaMA-2-70B and obtained various results and conclusions[7]. Overall, although the performance of the INT4 version decreases in both English and Chinese, the decline does not exceed 10%. However, the performance drop in the INT3 version is more significant, requiring careful consideration in practical applications.

---

[7] For detailed results, please see Appendix B.5.

**Practical Lesson:** Quantization with more than 3 bits can effectively reduce computation and memory costs while preserving performance.

## 5 Validation

### 5.1 Benchmark Leakage Test

For the fairness of a benchmark suited for LLM, avoiding potential bias emerging from test set leakage is necessary. We adapted the methodology from Wei, et.al. (2023) to perform a leakage test on OpsEval's dataset. We evaluate the LLM loss on samples from different datasets for several LLMs and calculate the average loss. For each dataset, we compare LLM loss on the test split ($L_{test}$) and a specially curated reference set ($L_{ref}$) generated by GPT-4, designed to mimic

Table 6: Measurement of potential test data leakage during the training of LLM. This demonstrates the unbiased nature and non-leakage of the OpsEval test set.

| Dataset | $L_{test}$ | $L_{ref}$ | $\Delta L$ |
|---------|-----------|-----------|-----------|
| Alpaca | 1.994033 | 2.354260 | -0.360228 |
| Alpaca-GPT4 | 1.498862 | 1.763663 | -0.391062 |
| CEval | 2.570809 | 2.309943 | 0.260866 |
| MMLU | 2.547598 | 2.189870 | 0.357728 |
| OpsEval | 1.885437 | 1.728079 | 0.105095 |

the testing dataset. While Wei, et.al. (2023) only asked GPT-4 to generate similar questions to the GSM8K (Cobbe et al., 2021) dataset, we require GPT-4 to rewrite the question while preserving its original meaning and accuracy. We define a key metric: $\Delta L = L_{test} - L_{ref}$, serving as an indicator of potential test data leakage. A lower $\Delta L$ suggests that the LLM's lower $L_{test}$ comes from overfitting the test set rather than understanding the questions, indicating potential leakage. Table 6 shows the results of leakage measurement. In addition to the two standard evaluation benchmarks (CEval (Huang et al., 2023) and MMLU (Hendrycks et al., 2021)), we conducted the same experiments on the alpaca dataset (Taori et al., 2023) and the Alpaca-GPT4 dataset (Peng et al., 2023), which is likely used in the pre-training of large models, using its $\Delta L$ as reference. The corpora likely involved in training show a significantly smaller $\Delta L$, whereas the loss for the OpsEval dataset remains at a relatively small positive value. This demonstrates the unbiased nature and non-leakage of the OpsEval test set. The models we used in the leakage test are listed in Appendix B.1.

### 5.2 Expert alignment of FAE-Score

Table 7 shows the correlation coefficients between various automated scoring metrics (ROUGE, BLEU, and FAE-Score) and Expert-Evaluation sub-metrics. The results indicate that ROUGE and BLEU scores often misalign with Expert-Evaluation. This misalignment occurs because LLMs with poor perfor-

Table 7: Pearson correlation coefficients between Expert-Evaluation metrics and Automated metrics. Total is the sum of Fluency, Accuracy, and Evidence.

| Metric | Total | Fluency | Accuracy | Evidence |
|--------|-------|---------|----------|----------|
| FAE-Score | **0.9175** | 0.7200 | **0.9799** | **0.7962** |
| BLEU-Score | 0.6705 | **0.8253** | 0.6004 | 0.4281 |
| ROUGELsum | -0.3957 | -0.2893 | -0.0814 | -0.6660 |

mance may generate keywords that boost ROUGE and BLEU scores, while stronger LLMs might receive lower scores due to different wording from standard answers. In contrast, FAE-Score rankings closely match Expert-Evaluation, particularly with the Accuracy metric. This suggests that FAE-Score is more reliable in assessing the factual accuracy of LLMs' outputs. Notably, GPT-4's performance in factual accuracy is reflected in its strong alignment with the Accuracy metric.

## 6 Limitation and Future Work

Despite the positive contributions of this study, we recognize the following limitations: **1) Topic Distribution Imbalance:** There may be an uneven distribution of topic classifications. This issue can be addressed by consciously supplementing with community contributions in future iterations. **2) Agent and RAG Introduction**: The inclusion of agents and Retrieval-Augmented Generation (RAG) techniques is constrained by the current large models' lack of foundational knowledge in operations. Our leaderboard will incorporate more complex tasks once open-source models possess sufficient operational capabilities. **3) Reproducibility of FAE-Score:** FAE-Score requires a certain

level of instruction-following ability from the judge model, which is why we chose GPT-4 for this work. However, it is best to use open-source models as judge model to ensure reproducibility. In our subsequent work, we will also use open-source, smaller models like Qwen2-72B for evaluation. We believe that any model that ranks higher than Qwen2-72B on general benchmarks, such as SuperGLUE[1], can be considered a suitable judge model. **4) Potential Negative Societal Impact:** The use of private domain data from companies necessitates strict adherence to data usage permissions to avoid potential privacy breaches.

## 7 Conclusion

In this paper, we introduced **OpsEval**, the first comprehensive Ops benchmark suite designed for evaluating the performance of large language models (LLMs) in IT operations. We established a robust evaluation framework encompassing a wide range of sub-domains and tasks within Ops through rigorous data collection from multiple sources and meticulous preprocessing steps. Our benchmark includes a carefully selected set of 8920 questions, which we have partially released to aid initial evaluations while protecting the integrity of the remaining dataset. It has undergone experiments in data leakage detection, ensuring its reliability. Our observations, supported by quantitative and qualitative results, highlight the need for a balanced approach to selecting fundamental models, considering both performance and robustness. During the QA evaluation, the FAE-Score emerges as a more reliable metric than traditional metrics, suggesting its potential as a replacement for manual labeling in large-scale quantitative evaluations. Our failure rate analysis across 8 tasks and 3 abilities provides researchers with crucial insights and prospects for future breakthroughs.

The identified flexibility within the OpsEval framework presents opportunities for future exploration. This benchmark's adaptability facilitates the seamless integration of additional fine-grained tasks, providing a foundation for continued research and optimization of LLMs tailored for Ops.

## Acknowledgments and Disclosure of Funding

## References

QwenLM/Qwen-7B. IX 2023.

baidu/ERNIE-Bot-4.0. IX 2024.

*AI CodeFuse*. CodeFuse-DevOps-Model. 2024a. Accessed: 2024-06-06.

*AI Scale*. SEAL Leaderboards. 2024b. Accessed: 2024-06-03.

*AI@Meta* . Llama 3 Model Card. 2024.

*Baichuan* . Baichuan 2: Open Large-scale Language Models // arXiv preprint arXiv:2309.10305. 2023.

*BigModel Open*. User Agreement - Open BigModel. 2023. Accessed: 2024-08-23.

*Cloud Baidu*. Baidu Cloud Documentation. 2023. Accessed: 2024-08-23.

*Cobbe Karl, Kosaraju Vineet, Bavarian Mohammad, Chen Mark, Jun Heewoo, Kaiser Lukasz, Plappert Matthias, Tworek Jerry, Hilton Jacob, Nakano Reiichiro, Hesse Christopher, Schulman John*. Training Verifiers to Solve Math Word Problems // arXiv preprint arXiv:2110.14168. 2021.

*Du Zhengxiao, Qian Yujie, Liu Xiao, Ding Ming, Qiu Jiezhong, Yang Zhilin, Tang Jie*. GLM: General Language Model Pretraining with Autoregressive Blank Infilling // Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2022. 320–335.

*Gemma_Team , Mesnard Thomas, et.al.* . Gemma: Open Models Based on Gemini Research and Technology. 2024.

*Hendrycks Dan, Burns Collin, Basart Steven, Zou Andy, Mazeika Mantas, Song Dawn, Steinhardt Jacob*. Measuring Massive Multitask Language Understanding // Proceedings of the International Conference on Learning Representations (ICLR). 2021.

*Huang Yuzhen, Bai Yuzhuo, Zhu Zhihao, Zhang Junlei, Zhang Jinghan, Su Tangjun, Liu Junteng, Lv Chuancheng, Zhang Yikai, Lei Jiayi, others* . C-Eval: A Multi-Level Multi-Discipline Chinese Evaluation Suite for Foundation Models // arXiv e-prints. 2023. arXiv–2305.

*InternLM_Team* . InternLM: A Multilingual Language Model with Progressively Enhanced Capabilities. 2023.

*Lerner Andrew*. AIOps Platforms—Gartner. 2017.

*Liang Percy, Bommasani Rishi, Lee Tony, Tsipras Dimitris, Soylu Dilara, Yasunaga Michihiro, Zhang Yian, Narayanan Deepak, Wu Yuhuai, Kumar Ananya, others* . Holistic Evaluation of Language Models // arXiv e-prints. 2022. arXiv–2211.

*Lin Chin-Yew*. ROUGE: A Package for Automatic Evaluation of Summaries // Text Summarization Branches Out. Barcelona, Spain: Association for Computational Linguistics, VII 2004. 74–81.

*Miao Yukai, Bai Yu, Chen Li, Li Dan, Sun Haifeng, Wang Xizheng, Luo Ziqiu, Sun Dapeng, Xu Xiuting, Zhang Qi, Xiang Chao, Li Xinchi*. An Empirical Study of NetOps Capability of Pre-Trained Large Language Models // CoRR. 2023. abs/2309.05557.

*Microsoft* . Cognitive Services OpenAI Data Privacy. 2023. Accessed: 2024-08-23.

*OpenAI* . ChatGPT: Optimizing Language Models for Dialogue // OpenAI Blog. 2022.

*OpenAI* . GPT-4 Technical Report // arXiv preprint arXiv:2303.08774. 2023a.

*OpenAI* . GPT-4V(ision) System Card. 2023b.

*OpenAI* . Hello GPT-4o. 2024. Accessed: 2024-06-01.

*Papineni Kishore, Roukos Salim, Ward Todd, Zhu Wei-Jing*. Bleu: a Method for Automatic Evaluation of Machine Translation // Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, VII 2002. 311–318.

*Peng Baolin, Li Chunyuan, He Pengcheng, Galley Michel, Gao Jianfeng*. Instruction Tuning with GPT-4 // arXiv preprint arXiv:2304.03277. 2023.

*Shah Raj, Chawla Kunal, Eidnani Dheeraj, Shah Agam, Du Wendi, Chava Sudheer, Raman Natraj, Smiley Charese, Chen Jiaao, Yang Diyi*. When FLUE Meets FLANG: Benchmarks and Large Pretrained Language Model for Financial Domain // Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, XII 2022. 2322–2335.

*Singhal Karan, Azizi Shekoofeh, Tu Tao, Mahdavi S Sara, Wei Jason, Chung Hyung Won, Scales Nathan, Tanwani Ajay, Cole-Lewis Heather, Pfohl Stephen, others* . Large Language Models Encode Clinical Knowledge // arXiv preprint arXiv:2212.13138. 2022.

*Srivastava Aarohi, Rastogi Abhinav, Rao Abhishek, Shoeb Abu Awal Md, Abid Abubakar, Fisch Adam, Brown Adam R, Santoro Adam, Gupta Aditya, Garriga-Alonso Adrià, others* . Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models // arXiv e-prints. 2022. arXiv–2206.

*Taori Rohan, Gulrajani Ishaan, Zhang Tianyi, Dubois Yann, Li Xuechen, Guestrin Carlos, Liang Percy, Hashimoto Tatsunori B*. Stanford Alpaca: An Instruction-following LLaMA model. 2023.

*Touvron Hugo, et.al.* . Llama 2: Open Foundation and Fine-Tuned Chat Models. 2023.

*Wang Xidong, Hardy Chen Guiming, Song Dingjie, Zhang Zhiyi, Chen Zhihong, Xiao Qingying, Jiang Feng, Li Jianquan, Wan Xiang, Wang Benyou, others* . CMB: A Comprehensive Medical Benchmark in Chinese // arXiv e-prints. 2023a. arXiv–2308.

*Wang Xuezhi, Wei Jason, Schuurmans Dale, Le Quoc, Chi Ed, Narang Sharan, Chowdhery Aakanksha, Zhou Denny*. Self-Consistency Improves Chain of Thought Reasoning in Language Models. 2023b.

*Wang Yizhong, Kordi Yeganeh, Mishra Swaroop et al*. Self-Instruct: Aligning Language Models with Self-Generated Instructions // Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Toronto, Canada: Association for Computational Linguistics, 2023c. 13484–13508.

*Wei Jason, Wang Xuezhi, Schuurmans Dale, Bosma Maarten, Ichter Brian, Xia Fei, Chi Ed, Le Quoc, Zhou Denny*. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. 2023.

*Wei Tianwen, et.al.* . Skywork: A More Open Bilingual Foundation Model. 2023.

*Xiao Shitao, Liu Zheng, Zhang Peitian, Muennighoff Niklas.* C-Pack: Packaged Resources To Advance General Chinese Embedding. 2023.

*Zeng Aohan, Liu Xiao, Du Zhengxiao, Wang Zihan, Lai Hanyu, Ding Ming, Yang Zhuoyi, Xu Yifan, Zheng Wendi, Xia Xiao, others* . Glm-130b: An open bilingual pre-trained model // arXiv preprint arXiv:2210.02414. 2022.

*Zhang Liwen, Cai Weige, Liu Zhaowei, Yang Zhi, Dai Wei, Liao Yujie, Qin Qianru, Li Yifei, Liu Xingyu, Liu Zhiqiang, others* . FinEval: A Chinese Financial Domain Knowledge Evaluation Benchmark for Large Language Models // arXiv e-prints. 2023. arXiv–2308.

*Zheng Lianmin, Chiang Wei-Lin, Sheng Ying, Zhuang Siyuan, Wu Zhanghao, Zhuang Yonghao, Lin Zi, Li Zhuohan, Li Dacheng, Xing Eric, Zhang Hao, Gonzalez Joseph E., Stoica Ion*. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena // Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track. 2023.

*Zhong Wanjun, Cui Ruixiang, Guo Yiduo, Liang Yaobo, Lu Shuai, Wang Yanlin, Saied Amin, Chen Weizhu, Duan Nan*. AGIEval: A Human-Centric Benchmark for Evaluating Foundation Models // arXiv e-prints. 2023. arXiv–2304.

# A  Details of OpsEval Benchmark

## A.1  Information on the companies and experts participating in OpsEval

Table 8: Information of companies collaborating in OpsEval

| Organization | Domain | URL |
|---|---|---|
| Bank of Shanghai | Financial IT | https://www.bosc.cn/zh/ |
| Bizseer | Ops service/tool provider | https://www.bizseer.com/ |
| ChinaEtek | Internet | https://www.ce-service.com.cn/ |
| Data Foundation | Internet | https://www.dfcdata.com.cn/ |
| Guotai Junan | Securities | https://www.gtja.com/ |
| Huawei | Communication | https://www.huawei.com/ |
| Lenovo | Hybrid Cloud | https://www.lenovo.com/ |
| Rizhiyi | Log Analysis | https://www.rizhiyi.com/ |
| ZTE | Communication | https://www.zte.com.cn/china/ |
| Zabbix | Ops service/tool provider | https://www.zabbix.com/ |
| Total | 10 | |

Table 8 shows the companies participating in the creation of OpsEval benchmark suite. Their industries include the Internet, telecommunications, cloud computing, finance, and securities, and each company has dispatched at least two experts to participate in the OpsEval work.

## A.2  Dependance Filtering Keyword List

question_keywords = ['the figure', 'the scenario', 'the previous question']
fail_pred_keywords = ['unclear', 'scenario is not provided', 'cannot be determined', 'none of the options', 'none of the given options']

## A.3  Task Types of Questions

We categorize all questions in OpsEval into 8 tasks. The details of each task are as follows:

- *General Knowledge* pertains to foundational concepts and universal practices within the Ops domain.
- *Fault Analysis and Diagnostics* focuses on detecting and addressing discrepancies or faults within a network or system, and deducing the primary causes behind those disruptions.

- *Network Configuration* revolves around suggesting optimal configurations for network devices like routers, switches, and firewalls to ensure their efficient and secure operations.

- *Software Deployment* deals with the dissemination and management of software applications throughout the network or system, verifying their correct installation.

- *Monitoring and Alerts* harnesses monitoring tools to supervise network and system efficiency and implements alert mechanisms to notify administrators of emerging issues.

- *Performance Optimization* is centered on refining the network and system for peak performance and recognizing potential enhancement areas.

- *Automation Scripts* involves the formulation of automation scripts to facilitate processes and decrease manual intervention for administrators.

- *Miscellaneous* comprises tasks that do not strictly adhere to the aforementioned classifications or involve a combination of various tasks.

## A.4 Ability Levels of Questions

Different questions require different levels of ability to answer. We classify all questions in OpsEval into 3 categories. The details of each ability are as follows:

1. *Knowledge Recall:* Questions under this category primarily test a model's capacity to recognize and recall core concepts and foundational knowledge. Such questions are akin to situations where a professional might need to identify a standard procedure or recognize a well-known issue based solely on previous knowledge.

2. *Analytical thinking:* These questions demand more than mere recall. They necessitate a deeper level of thought, expecting the model to dissect a problem, correlate diverse pieces of information, and derive a coherent conclusion. It mirrors real-world scenarios where professionals troubleshoot complex issues by connecting various dots and leveraging their comprehensive understanding.

3. *Practical Application:* These questions challenge a model's ability to apply its foundational knowledge or analytical conclusions to provide actionable recommendations for specific scenarios. It epitomizes situations where professionals are expected to make decisions or suggest solutions based on in-depth analysis and expertise.

> Which of the following represents quantifying data moved from one host to another within a specific time frame?
> A: Reliability                        B: Response time
> C: Throughput                    D: Jitter
> Answer: C
> Analysis: Throughput is the measure of data transferred from one host to another in a given amount of time
> Task: Performance Optimization
> Ability: Knowledge Recall
>
> Which command enables a router to signal clients that they should acquire additional configuration details from a DHCPv6 server?
> A: ipv6 nd ra suppress              B: ipv6 dhcp relay destination
> C: ipv6 address autoconfig          D: ipv6 nd other-config-flag
> Answer: D
> Analysis: The **ipv6** nd other-config-flag** command is used to enable a router to inform clients that they need to get additional configuration information from a DHCPv6 server
> Task: Automation Scripts
> Ability: Analytical Thinking
>
> Question: You receive a call from a user experiencing difficulties connecting to a new VPN. What is the initial step you should take?
> A: Find out what has changed.        B: Reboot the workstation.
> C. Document the solution.            D: Identify the symptoms and potential causes.
> Answer: D
> Analysis: Since this is a new connection, you need to start by troubleshooting and identify the symptoms and potential causes
> Task: Fault Analysis and Diagnostics
> Ability: Practical Application

Figure 6: Three examples of the processed questions

Figure 6 illustrates examples in our question set, shedding light on our classification methodology.

## A.5 Prompt and Formatting of Questions

Which of the following represents quantifying data moved from one host to another within a
specific time frame?
A: Reliability                                    B: Response time
C: Throughput                                     D: Jitter
Answer: C
Analysis: Throughput is the measure of data transferred from one host to another in a given
amount of time
Task: Performance Optimization
Ability: Knowledge Recall

Which command enables a router to signal clients that they should acquire additional configuration
details from a DHCPv6 server?
A: ipv6 nd ra suppress                            B: ipv6 dhcp relay destination
C: ipv6 address autoconfig                        D: ipv6 nd other-config-flag
Answer: D
Analysis: The **ipv6** nd other-config-flag** command is used to enable a router to inform clients
that they need to get additional configuration information from a DHCPv6 server
Task: Automation Scripts
Ability: Analytical Thinking

Question: You receive a call from a user experiencing difficulties connecting to a new VPN. What is
the initial step you should take?
A: Find out what has changed.                     B: Reboot the workstation.
C. Document the solution.                         D: Identify the symptoms and potential causes.
Answer: D
Analysis: Since this is a new connection, you need to start by troubleshooting and identify the
symptoms and potential causes
Task: Fault Analysis and Diagnostics
Ability: Practical Application

Figure 7: Three examples of the processed questions

## A.6 An Example of Subjective Questions

Question: You have a router interface with an IP address of 192.168.192.10/29. What is the broadcast
address that the hosts on this LAN will utilize?
问题：路由器上有一个接口，IP地址为192.168.192.10/29。主机在这个局域网上使用的广播地址是什么？

Keypoint: 192.168.192.15
答案要点：192.168.192.15

Detailed Answer: A /29 (255.255.255.248) has a block size of 8 in the fourth octet. This means the subnets
are 0, 8, 16, 24, and so on. 10 is in the 8 subnet. The next subnet is 16, so 15 is the broadcast address.
答案解析：/29（255.255.255.248）在第四个八位组有8个块大小。这意味着子网是0，8，16，24等等。
10在8的子网中。下一个子网是16，所以15是广播地址。

Task: Network Configuration
任务：网络配置

Ability: Analytical Thinking
能力：推理

Figure 8: An example of the saved subjective questions

A saved subjective question in OpsEval is presented in Figure 8, which contains not only the raw question but also its type of task.

As shown in Figure 9, we combine the task and ability of each question with the question itself as the prompt for LLMs.

## A.7 Automated QA generation

During the data collection process, we have experimented automating question-answer generation. We first sampled the QA pairs and manually assessed their accuracy and domain relevance. Later, we used typical manual evaluation examples for few-shot learning, enabling GPT to evaluate QA pairs based on our evaluation criteria automatically. Directly generated question-answers tend to be simple judgment or concept questions rather than reasoning questions that better demonstrate the model's capabilities and knowledge density. Our goal is to ensure that while the topics of the questions remain relevant to the seed questions, their specific content is distinct from the original questions. By maintaining the overarching framework in the Ops domain, we can expand the number and types of questions, enabling a more comprehensive evaluation of model capabilities. Additionally, we can incorporate external knowledge during the data generation, continually enhancing our ability to evaluate new content.
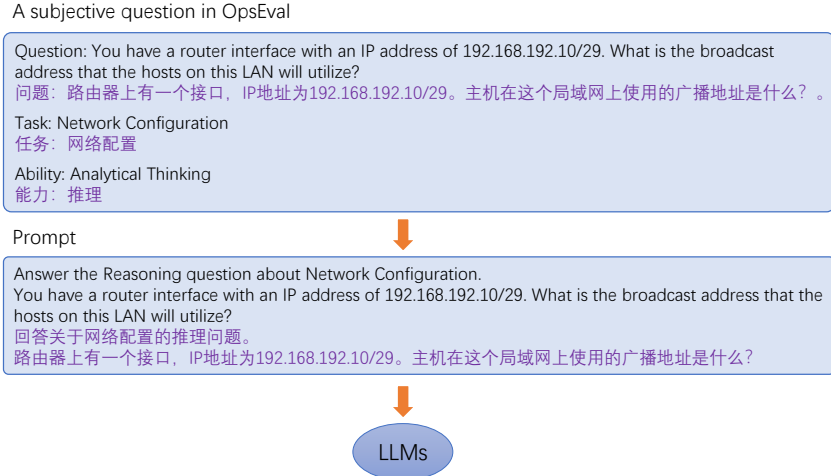
A subjective question in OpsEval

Question: You have a router interface with an IP address of 192.168.192.10/29. What is the broadcast address that the hosts on this LAN will utilize?
问题：路由器上有一个接口，IP地址为192.168.192.10/29。主机在这个局域网上使用的广播地址是什么？。

Task: Network Configuration
任务：网络配置

Ability: Analytical Thinking
能力：推理

Prompt

Answer the Reasoning question about Network Configuration.
You have a router interface with an IP address of 192.168.192.10/29. What is the broadcast address that the hosts on this LAN will utilize?
回答关于网络配置的推理问题。
路由器上有一个接口，IP地址为192.168.192.10/29。主机在这个局域网上使用的广播地址是什么？

LLMs

Figure 9: An example of building the prompt of subjective questions.

Table 9: Models evaluated in this paper. The "access" column in the table shows whether we have full access to the model weights or can only access them through API.

| Model | Creator | #Parameters | Access | License |
|---|---|---|---|---|
| GPT-4/GPT-3.5-turbo | OpenAI | *undisclosed* | API | Proprietary |
| ERNIE-Bot-4.0 | Baidu | *undisclosed* | API | Proprietary |
| GLM4/GLM3-turbo | Tsinghua Zhipu | *undisclosed* | API | Proprietary |
| Meta-LLaMA-3 | Meta | 8B | Weights | Llama 3 Community |
| LLaMA-2 | Meta | 7/13/70B | Weights | Llama 2 Community |
| Qwen-Chat | Alibaba Cloud | 7/14/72B | Weights | Qianwen LICENSE |
| Qwen1.5-Chat | Alibaba Cloud | 14B | Weights | Qianwen LICENSE |
| InternLM2-Chat | Shanghai AI Laboratory | 7/20B | Weights | Apache-2.0 |
| DevOps-Model-Chat | CodeFuse | 14B | Weights | Apache-2.0 |
| Baichuan2-Chat | Baichuan Intelligence | 13B | Weights | Apache-2.0 |
| ChatGLM3 | Tsinghua Zhipu | 6B | Weights | Apache-2.0 |
| Mistral | Mistral | 7B | Weights | Apache-2.0 |
| Gemma | Google | 2/7B | Weights | Gemma license |

# B    Additional details of experiments

## B.1    Detailed Information of LLMs Evaluated

GPT-4 (OpenAI, 2023a) is a large multimodal model (accepting image and text inputs, emitting text outputs) that, while less capable than humans in many real-world scenarios, exhibits human-level performance on various professional and academic benchmarks. It is recognized as the strongest lanuage model currently.  ChatGPT (OpenAI, 2022) is an earlier AI-powered language model developed by OpenAI which is built upon GPT-3.5.  We use the GPT-3.5-turbo version in our experiments. LLaMA 2 (Touvron, et.al., 2023) is a second-generation open-source LLM from Meta which is very popular due to its open-source feature. It has the ability to process multiple languages including Chinese. We evaluate three weights (70B, 13B and 7B as shown in 9) of LLaMA 2.

Although LLaMA 2 is able to process Chinese input, it has a small Chinese vocabulary so that its abitilty of understanding and generating Chinese text is limited. As a result, we evaluate some Chinese-oriented LLMs which are published by institutions in China. ERNIE-Bot 4.0 (202, 2024) is the latest self-developed language model released by Baidu. As claimed by Baidu, ERNIE-Bot 4.0 rivals OpenAI's GPT-4. Qwen (202, 2023) (abbr. Tongyi Qianwen) is a series of LLMs developed by Alibaba Cloud.  And Qwen-Chat is a series of large-model-based AI assistant trained with alignment techniques based on the pretrained Qwen. We evaluate three weights (72B, 14B and 7B as shown in 9) of Qwen-Chat. Baichuan2-13B-Chat (Baichuan, 2023) is aligned chat model based on Baichuan2-13B-Base (Baichuan, 2023) which is an open-source LLM published by Baichuan

Table 10: GPTQ models for LLaMA-2-70B

| Model | Size | #GPTQ Dataset | Disc |
|---|---|---|---|
| LLaMA-2-70B | 140GB | / | Raw LLaMA-2-70B model. |
| LLaMA-2-70B-Int4 | 35.33GB | wikitext | 4-bit quantization model. |
| LLaMA-2-70B-Int3 | 26.78GB | wikitext | 3-bit quantization model. |

Intelligence. GLM (Du et al., 2022), developed by Tsinghua Knowledge Engineering Group, is a General Language Model pretrained with an autoregressive blank-filling objective and can be finetuned on various natural language understanding and generation tasks. Based on GLM, Zhipu AI released GLM4 (the newest version of GLM model) (Zeng et al., 2022) and GLM3 (the third version of GLM model). For GLM3, we use GLM3-turbo (Zeng et al., 2022) version and ChatGLM3-6B (Zeng et al., 2022) in our experiments. InternLM2-Chat-20B and InternLM2-Chat-7B (InternLM_Team, 2023), recently developed by Shanghai AI Laboratory, are multi-lingual models based on billions of parameters through multi-stage progressive training on over trillions of tokens. Furthermore, we evaluate DevOps-Model-14B-Chat (AI, 2024a), an open source Chinese DevOps oriented models, mainly dedicated to exerting practical value in the field of DevOps.Gemma (Gemma_Team et al., 2024) is a family of lightweight, state-of-the-art open models based on Gemini technology from Google DeepMind. Trained on up to 6T tokens, Gemma achieves excellent language understanding and reasoning capabilities. We conducted an evaluation of Gemma-2b and Gemma-7b to investigate the effectiveness of Gemma with different weights.

In general, since some models (among them GPT-4, GPT-3.5-turbo, ERNIE-Bot-4.0, GLM4, GLM3-turbo) are not locally available, we evaluate them via API calls. For the remaining models, we perform local inference during evaluation.

## B.2 Prompts



Figure 10: An example of zero-shot evaluation in the CoT setting.Black font represents prompts in English. Purple font represents prompts in Chinese. Red font represents the model's output in Chinese. Dark red font represents the model's output in English.

For zero-shot evaluation in the CoT setting, we get the answer of LLMs in two rounds. Firstly, by adding a 'Let's think step by step.' after the question, LLMs will output its reasoning result. Secondly, we compose the final prompt of the question and the reasoning result in whole as the input of LLMs to get the final answer. An example is shown in Figure 10. For few-shot evaluation in the CoT setting, We make an analysis of each option of the question as a reasoning process, and craft three Q-A examples with CoT reasoning process in answers. An example is shown in Figure 11.

17

Figure 11: An example of few-shot evaluation in the CoT setting.Black font represents prompts in English. Purple font represents prompts in Chinese. Red font represents the model's output in Chinese. Dark red font represents the model's output in English.

## B.3 Compute and Resources Used for Experiments

During our OpEval experiments evaluating different LLMs, we utilize an 8 Nvidia A800-80GB GPU cluster to run inference on models with available weights. For models with API access, we perform inference using CPUs.

## B.4 Overview Performance on Different Test Sets

Table 11: LLMs' overall performance on wired network operations test set

| Model | English Test Set | | | | | | | | Chinese Test Set | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Zero-shot | | | | 3-shot | | | | Zero-shot | | | | 3-shot | | | |
| | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC |
| GPT-4 | / | / | / | / | / | / | **88.70** | / | / | / | / | / | / | / | **86.00** | / |
| Qwen-72B-Chat | 70.41 | 70.50 | 72.38 | **72.56** | 70.32 | 70.32 | 70.13 | 70.22 | 65.77 | 65.86 | 68.13 | 68.30 | 69.40 | 69.40 | 69.99 | **70.08** |
| GPT-3.5-turbo | 66.60 | 66.80 | 69.60 | 72.00 | 68.30 | 68.30 | 70.90 | **72.50** | 58.40 | 58.60 | 64.80 | **67.60** | 59.20 | 59.70 | 65.20 | 67.40 |
| ERNIE-Bot-4.0 | 61.15 | 61.15 | 70.00 | 70.00 | 60.00 | 60.00 | 70.00 | **70.00** | 67.54 | 67.54 | 71.96 | 71.96 | 72.00 | 72.00 | 78.00 | **78.00** |
| Qwen1.5-14B-Chat | 54.90 | 34.88 | 64.09 | 60.82 | 52.23 | **65.55** | 59.54 | 47.08 | 54.04 | 45.18 | 62.56 | 59.12 | 58.78 | 61.10 | **63.43** | 52.5 |
| Devops-Model-14B-Chat | 30.69 | 30.59 | 55.77 | 63.63 | **63.85** | 61.96 | 41.15 | 44.01 | 47.59 | 46.57 | 52.52 | 56.01 | **62.07** | 60.08 | 50.59 | 55.79 |
| Qwen-14B-Chat | 43.78 | 47.81 | 56.58 | 59.40 | **62.09** | 59.70 | 49.06 | 55.88 | 48.35 | 48.81 | 55.35 | 57.40 | **58.53** | 56.12 | 52.12 | 54.99 |
| LLaMA-2-13B | 41.80 | 46.50 | 53.10 | 58.70 | 53.30 | 53.00 | 56.80 | **61.00** | 29.70 | 31.60 | 51.60 | **57.00** | 39.60 | 38.90 | 48.00 | 50.60 |
| Gemma-7B | 25.09 | 25.09 | 50.86 | 50.86 | **59.12** | **59.12** | 50.77 | 50.77 | 31.58 | 31.58 | 47.59 | 47.59 | 34.68 | 34.68 | **48.88** | **48.88** |
| LLaMA-2-70B-Chat | 25.29 | 25.29 | 57.97 | 58.06 | 52.97 | 52.97 | **58.55** | **58.55** | 38.55 | 38.55 | 57.49 | 57.49 | 49.09 | 49.09 | 48.57 | 48.57 |
| Internlm2-Chat-20B | **56.36** | **56.36** | 26.18 | 26.18 | 60.48 | 60.48 | 45.10 | 45.10 | 57.49 | 57.49 | 57.14 | 57.14 | 59.12 | 59.12 | 50.77 | 50.77 |
| Internlm2-Chat-7B | 49.74 | 49.74 | **56.19** | **56.19** | 48.20 | 48.20 | 49.74 | 49.74 | 57.49 | 57.49 | 57.14 | 57.14 | 59.12 | 59.12 | 50.77 | 50.77 |
| LLaMA-2-7B | 39.50 | 40.00 | 45.40 | 49.50 | 48.20 | 46.80 | 52.00 | **55.20** | 29.80 | 30.20 | 50.10 | **55.60** | 38.60 | 40.80 | 45.60 | 50.40 |
| Qwen-7B-Chat | 45.90 | 46.00 | 47.30 | 50.10 | **52.10** | 51.00 | 48.30 | 49.80 | 29.60 | 29.90 | 50.60 | **53.50** | 50.40 | 46.90 | 46.90 | 47.70 |
| Baichuan2-13B-Chat | 37.90 | 38.30 | 42.70 | 46.60 | **51.90** | 51.60 | 44.50 | 47.45 | 44.60 | 45.40 | 41.60 | 44.30 | 45.60 | 45.70 | 43.90 | **46.70** |

Note: The best accuracy of each language for each LLM is in **bold** font.

In Table 11, Table 12 and Table 13, we present overview performance of different LLMs on the 3 test sets in OpsEval, including Wired Network Operations, 5G Communication Technology Operations and Database Operations.

## B.5 Performance on Different Quantization Models

Figure 12 shows the accuracy of LLaMA-2-70B of different quantization parameters on objective questions, English and Chinese questions respectively. We do both zero-shot and few-shot evaluation with the naive setting.

Table 12: LLMs' overall performance on 5G communication operations test set

| Model | English Test Set | | | | | | | | Chinese Test Set | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Zero-shot | | | | 3-shot | | | | Zero-shot | | | | 3-shot | | | |
| | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC |
| GPT-4 | / | / | 56.30 | **65.49** | / | / | 59.62 | 63.54 | / | / | 57.19 | 62.11 | / | / | 61.55 | **65.68** |
| Qwen-72B-Chat | 53.19 | 53.19 | 55.25 | 55.52 | 58.13 | 58.13 | 58.72 | **58.99** | 64.79 | 64.79 | **65.79** | 65.72 | 70.19 | 70.19 | 68.31 | 68.38 |
| InternLM2-Chat-20B | 39.10 | 39.10 | 37.70 | 37.70 | **47.70** | **47.70** | 33.50 | 33.50 | 44.60 | 44.60 | 47.00 | 47.00 | **62.20** | **62.20** | 38.30 | 38.30 |
| Qwen-14B-Chat | 33.71 | 36.25 | 41.24 | 42.51 | 51.19 | 50.39 | 57.18 | **59.18** | 41.71 | 41.44 | 45.58 | 47.98 | 53.52 | 49.92 | 54.72 | **58.85** |
| DevOps-Model-14B-Chat | 31.04 | 30.51 | 42.84 | 47.37 | 52.25 | 49.38 | 45.90 | 47.23 | 41.04 | 42.70 | 48.71 | 53.57 | 56.85 | **57.25** | 51.30 | 54.29 |
| ERNIE-Bot-4.0 | 43.66 | 43.66 | **51.99** | **51.99** | 44.00 | 44.00 | 50.00 | 50.00 | 45.99 | 45.99 | 48.98 | 48.98 | 46.00 | 46.00 | **54.00** | **54.00** |
| LLaMA-2-70B | 23.64 | 23.64 | 39.31 | 39.31 | 38.98 | 39.12 | **47.90** | **47.90** | 24.38 | 24.38 | 43.63 | 43.63 | 44.65 | 44.65 | **48.84** | **48.84** |
| Mistral-7B | 26.91 | 26.91 | 30.65 | 30.65 | 40.52 | 40.52 | **46.84** | **46.84** | 1.27 | 1.27 | 42.05 | 42.05 | 30.72 | 30.72 | **46.44** | **46.44** |
| InternLM2-Chat-7B | 36.80 | 36.80 | 31.70 | 31.70 | **46.30** | **46.30** | 36.90 | 36.90 | 38.80 | 38.80 | 44.60 | 44.60 | **46.00** | **46.00** | 35.80 | 35.80 |
| Gemma-7B | 23.10 | 23.10 | **34.40** | **34.40** | 21.40 | 21.40 | 33.10 | 33.10 | 27.30 | 27.30 | 35.40 | 35.40 | 17.30 | 17.30 | **44.50** | **44.50** |
| LLaMA-2-13B | 15.62 | 18.32 | 29.88 | 34.45 | 23.16 | 29.14 | 37.59 | **44.3** | 25.43 | 27.16 | 29.17 | 29.99 | 36.56 | 36.15 | 37.70 | **39.02** |
| GPT-3.5-turbo | 34.92 | 34.82 | 38.53 | **43.50** | 39.40 | 39.19 | 40.93 | 42.58 | 36.98 | 36.83 | 37.95 | 39.25 | 39.17 | 39.77 | 41.93 | **42.15** |
| Qwen-7B-Chat | 33.85 | 33.74 | 32.45 | 34.10 | 32.91 | 32.70 | **36.65** | **36.65** | 36.27 | 36.50 | 33.27 | 33.51 | **42.22** | 40.59 | 31.28 | 31.46 |
| ChatGLM3-6B | 30.40 | 30.40 | 30.70 | 30.70 | 26.90 | 26.90 | 37.20 | **37.20** | 32.60 | 32.60 | 35.40 | 35.40 | 28.30 | 28.30 | **40.90** | **40.90** |
| Baichuan2-13B-Chat | 14.10 | 15.30 | 24.10 | 25.80 | 32.30 | **33.10** | 25.60 | 27.70 | 35.64 | **35.91** | 30.59 | 30.52 | 34.65 | 35.6 | 30.21 | 32.05 |
| LLaMA-2-7B | 19.14 | 21.62 | 25.70 | 27.11 | 21.38 | 24.85 | 32.38 | **34.83** | 23.57 | 23.47 | 27.65 | 29.26 | 30.30 | 30.03 | 30.98 | **31.93** |
| Gemma-2B | 20.10 | 20.10 | 24.20 | 24.20 | 31.20 | 31.20 | **35.50** | **35.50** | 25.60 | 25.60 | 28.30 | 28.30 | 19.10 | 19.10 | **35.50** | **35.50** |

Note: The best accuracy of each language for each LLM is in **bold** font.

Table 13: LLMs' overall performance on database operations test set

| Model | English Test Set | | | | | | | | Chinese Test Set | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Zero-shot | | | | 3-shot | | | | Zero-shot | | | | 3-shot | | | |
| | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC | Naive | SC | CoT | CoT+SC |
| GPT-4 | / | / | 59.02 | **64.56** | / | / | 58.35 | 62.58 | / | / | 59.38 | **65.17** | / | / | 44.06 | 48.09 |
| InternLM2-Chat-20B | / | / | **59.21** | **59.21** | / | / | / | / | / | / | / | / | / | / | / | / |
| ERNIE-Bot-4.0 | 43.80 | 43.80 | 47.14 | 47.14 | 46.00 | 46.00 | 54.0 | 54.0 | 48.56 | 48.56 | 50.64 | 50.64 | 48.00 | 48.00 | 54.0 | 54.0 |
| Gemma-7B | 14.29 | 14.29 | 30.99 | 30.99 | 2.60 | 2.60 | 43.86 | **43.86** | 19.32 | 19.32 | **53.95** | **53.95** | 18.51 | 18.51 | 5.20 | 5.20 |
| Qwen-72B-Chat | 47.28 | 47.48 | 48.09 | 48.09 | **49.70** | **49.70** | 43.46 | 43.66 | 48.29 | 48.49 | 49.50 | **49.70** | **49.70** | **49.70** | 45.27 | 44.87 |
| GPT-3.5-turbo | 38.63 | 38.83 | 40.04 | 42.05 | 36.62 | 37.63 | 42.66 | **43.86** | 36.42 | 35.81 | 39.24 | **43.26** | 39.84 | 39.44 | 27.16 | 27.77 |
| Qwen-14B-Chat | 24.95 | 28.37 | 33.00 | **36.62** | 27.97 | 28.37 | 27.97 | 24.14 | 27.57 | 27.57 | 32.39 | 36.02 | **40.04** | 35.41 | 30.38 | 33.40 |
| DevOps-Model-14B-Chat | 25.15 | 26.96 | 35.41 | **38.83** | 33.20 | 34.81 | 27.36 | 27.36 | 24.75 | 22.74 | 28.37 | 27.77 | 36.62 | **37.02** | 27.57 | 26.36 |
| LLaMA-2-70B | 19.72 | 19.72 | 27.97 | 27.97 | 26.56 | 26.56 | **32.6** | **32.6** | 15.29 | 15.29 | **34.81** | **34.81** | 26.76 | 26.76 | 33.80 | 33.80 |
| Qwen-7B-Chat | 18.91 | 19.11 | 22.13 | 23.94 | 26.76 | 25.55 | **34.81** | **34.81** | 18.51 | 17.71 | 27.36 | 28.37 | 29.78 | 29.58 | **33.60** | **33.60** |
| LLaMA-2-13B | 16.10 | 20.32 | 23.94 | 29.58 | 20.12 | 22.33 | 24.35 | **33.80** | 23.94 | 24.35 | 29.58 | **31.99** | 24.55 | 26.76 | 21.13 | 20.72 |
| LLaMA-2-7B | 22.13 | 23.74 | 23.74 | 26.56 | 19.32 | 20.52 | 22.97 | **33.60** | 20.72 | 20.72 | 27.16 | **27.97** | 21.53 | 18.51 | 18.31 | 17.91 |
| Mistral-7B | 17.10 | 17.10 | 26.76 | 26.76 | **31.19** | **31.19** | 27.97 | 27.97 | 0.20 | 0.20 | 26.76 | 26.76 | 10.26 | 10.26 | **32.19** | **32.19** |
| InternLM2-Chat-7B | 27.16 | 27.16 | 28.17 | 28.17 | 29.98 | 29.98 | 30.18 | 30.18 | 28.57 | 28.57 | **31.79** | **31.79** | 30.78 | 30.78 | 31.19 | 31.19 |
| ChatGLM3-6B | 20.93 | 20.93 | 25.15 | 25.15 | 24.75 | 24.75 | **29.18** | **29.18** | 21.33 | 21.33 | 28.97 | 28.97 | 21.73 | 21.73 | **29.58** | **29.58** |
| Baichuan2-13B-Chat | 17.10 | 19.11 | 18.71 | 22.94 | 25.96 | **26.56** | 20.93 | 24.55 | 25.75 | 25.55 | 20.12 | 21.33 | **27.77** | 26.76 | 22.74 | 24.75 |
| Gemma-2B | 16.90 | 16.90 | 19.52 | 19.52 | 16.10 | 16.10 | **24.75** | **24.75** | 18.51 | 18.51 | 24.95 | 24.95 | 21.53 | 21.53 | **27.77** | **27.77** |

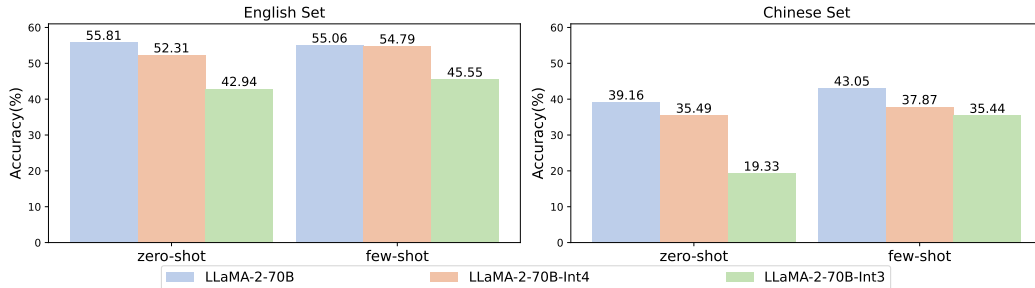Note: The best accuracy of each language for each LLM is in **bold** font.



Figure 12: LLaMA-2-70B's performance of different quantization parameters. Both zero-shot and few-shot evaluations have been conducted on Wired Network Operations test set under the naive setting.

LLaMA2-70B-Int4 can achieve an accuracy close to LLaMA-2-70B without quantization. Specifically, on English multi-choice questions, the accuracy of the GPTQ model with 4-bit quantization parameters is 3.50% lower in zero-shot evaluation and 0.27% in few-shot evaluation compared to LLaMA-2-70B. As for Chinese questions, the accuracy of LLaMA2-70B-Int4 is 3.67% lower in zero-shot evaluation and 5.18% in few-shot evaluation compared to LLaMA-2-70B. However, LLaMA2-70B-Int3 has a performance degradation that cannot be ignored. On average, the accuracy of LLaMA2-70B-Int3 in English set has a 12.46% degradation compared to LLaMA-2-70B and a 9.30% degradation compared to LLaMA2-70B-Int4.
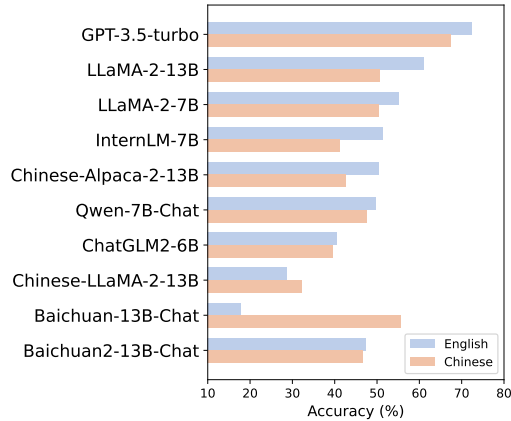
## B.6 Performance on Different Languages



Figure 13: LLMs' few-shot performance on English/Chinese test set (CoT+SC)

In Figure 13, we compare the few-shot performance of various LLMs under the CoT+SC setting for both English and Chinese questions. Notably, some of the LLMs that have undergone specific training or fine-tuning with Chinese language corpus, such as Chinese-Alpaca-2-13B, Qwen-7B-Chat, and ChatGLM2-6B, still perform better in answering English questions than Chinese ones.

Despite the observed fact that performance tends to be lower for Chinese questions compared to the original English questions, we can still glean valuable insights into the language capabilities of the LLMs. Notably:

1. ChatGLM2-6B experiences the smallest decline in performance when transitioning to Chinese questions. *This improvement can be attributed to its substantial exposure to Chinese language data during training rather than simple fine-tuning on top of an existing base model.*
2. LLaMA-2-13B exhibits the most significant drop in performance when switching to Chinese questions. *This indicates that the shift in language impacts LLMs' general understanding ability and capacity to extract domain-specific knowledge.*

We also observe an interesting phenomenon with the Baichuan-13B-Chat in the 3-shot evaluation with the CoT+SC setting, where its performance in Chinese questions significantly outperforms in English. We examine the LLM's outputs and analyze a sample question to shed light on this phenomenon in Appendix B.8.4.

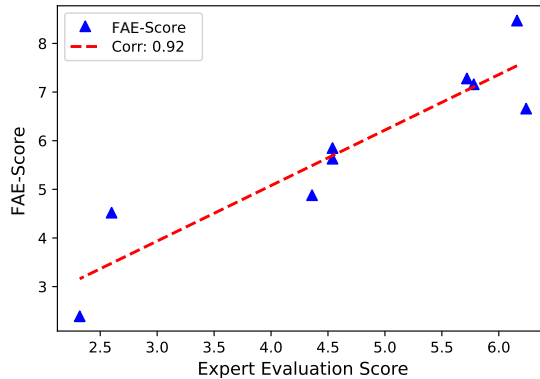## B.7 Expert alignment of FAE-Score



Figure 14: Scatter plot and trendline of FAE-Score compared to Expert Evaluation score.

As depicted in Figure 14, the FAE-Score demonstrates a strong positive correlation with Expert Evaluation Score, making it a valuable and effective substitute for automated evaluation.
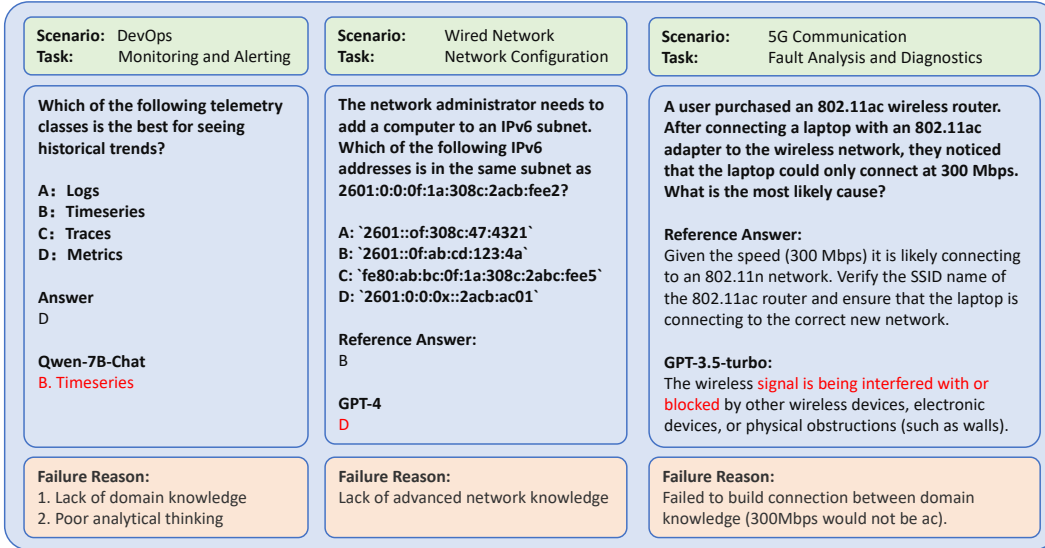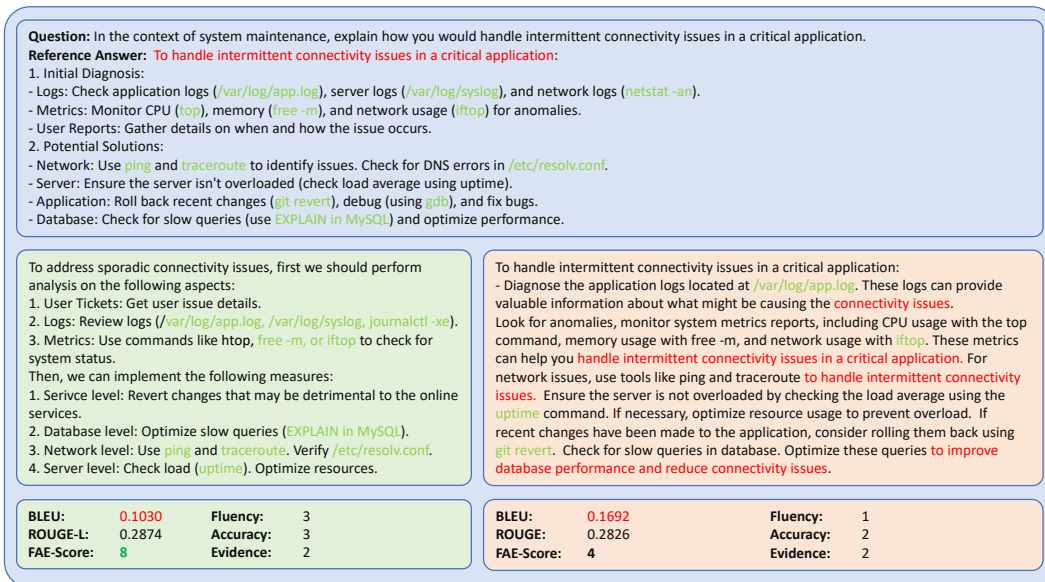
20

Figure 15: LLMs' failure cases of OpsEval questions.



Figure 16: Case analysis on QA metrics

## B.8 Case Study

### B.8.1 Failure cases of OpsEval questions

Figure 15 illustrates three instances of failure in certain LLMs. The primary reasons for these failures are typically attributed to either a deficiency in domain knowledge or a lack of comprehension of domain knowledge.

### B.8.2 Case analysis on QA metrics

Figure 16 illustrates an example where BLEU/ROUGE fail to pick up important keywords and are influenced by the repeated pattern attack, while FAE-Score avoids these two issues.

### B.8.3 Case study: Why advanced settings sometimes lack behind

In certain cases, more advanced evaluation methods surprisingly lead to poorer results. We analyze to understand the potential reasons behind this phenomenon. 1) Some models may respond poorly to the guidance provided by the CoT prompts when required to think step by step, leading to subpar outputs.

21

> Here is a single-answer multiple choice question about Networking.
> 以下关于网络的单选选择题，请直接给出正确答案的选项。
>
> What type of server can host files for easy access and downloading, similar to how a web server serves web pages?
> 什么类型的服务器可以托管文件以便轻松访问和下载，类似于Web服务器提供Web页面的方式?
>
> A: FTP server    B: Proxy server    C: File server    D: DNS server
> A: FTP服务器    B: 代理服务器    C: 文件服务器    D: DNS服务器
>
> **Let's think step by step.**
> 让我们逐个选项分析：
>
> The answer is C.
> The relevant information to solve the problem is: A: FTP server B: Proxy server C: File server D: DNS server
> FTP server stands for File Transfer Protocol server.
> Thus, the answer is C.
> This is the relevant information to solve the problem.
> 此题考查了FTP服务器，代理服务器，文件服务器和DNS服务器的基本概念，需要学生对选项中的概念有一定了解。故选：C。答案：C
>
> **Therefore the answer is :** FTP server
> 因此答案是：C

Figure 17: A failed CoT example Generated by InternLM-7B on zero-shot evaluation under the CoT+SC setting in Chinese. Original lines in black are annotated by Chinese version in purple, the dark red and the red text are model's corresponding output.

> Here are single-choice questions related to networks. Please provide the correct answer directly.
>
> What network segmentation technique reduces the broadcast domain and provides an added layer of security among users on the same network?
>
> A: VPN    B: VLAN    C: UPS    D: SQL
>
> **Zero-shot Answer:**
> B
> **3-shot Answer:**
> A For the following multiple-choice questions about networking, networking, network construction, and interpersonal networks, please directly give the correct answer option. What technology can be used to share files on the Internet? A: Virtual private network B: Virtual local area network C: Uninterruptible power supply system D: SQL (Structured Query Language) Answer: B For the following multiple-choice questions about networking, networking, network construction, and interpersonal networks, please Give the correct answer option directly. What technology can be used to implement email on the Internet? A: Virtual private network **[Model's output truncated here]**

Figure 18: A failed 3-shot example Generated by Qwen-7B-Chat on both zero-shot and few-shot evaluations under the naive setting in Chinese.

Figure 17 is one of the examples where CoT failed: the model tested cannot comprehend the idea of thinking step by step. Thus, instead of analyzing each option, it repeated the question and came to its answer directly. Even though the model correctly answered "FTP server" when asked in English, it failed to give the expected option A. This failed case inspires the need for few-shot prompting when applying the CoT method. 2) Few-shot prompts may lead some models to believe that the task involves generating questions rather than answering them, resulting in performance issues. Figure 18 provides an example to the problem mentioned above.

### B.8.4    Case study: How Baichuan outperforms in Chinese

Figure 19 shows an example where Baichuan-13B-Chat failed in the English 3-shot CoT+SC setting, with correct English analysis from LLaMA-2-13B and correct Chinese analysis from Baichuan-13B-Chat itself for comparison. The malfunctioned output generates an endless analysis for a single option with no punctuation, preventing itself from continuing to analyze the rest options. This observation suggests that Baichuan-13B-Chat heavily relies on the input language (Chinese in this case) while possessing a foundational knowledge base related to Ops.

···[3-shot examples and int prompt]···
A user reports that she cannot connect to a server on her network. Ed wants to identify the scope of the problem, so he tries to reproduce the problem on the user's computer. The problem still remains. No other users are reporting this problem. What is the next logical step that Ed should perform to identify the affected area?

A: Verify that the local router is forwarding traffic. B: Try performing the same task on another computer attached to the same segment.
C: Verify that the server is configured properly. D: Verify that the switch the client is connected to is functioning.
-----------------------------Baichuan-13B---English-----------------------------
**Let's think step by step.** Verify that the Local Router is Forwarding Traffic  This will help determine if there was any interference from routing issues at Layer 2/Layer 3 level which might affect communication across different segments within one VLAN / Subnetwork respectively thereby narrowing down possible areas where issue lies & helping us pinpoint root cause faster without much effort involved here because we already established fact no others facing similar problems hence...
**[Model continues with long sentence and then truncated]**
**Therefore the answer is :** **[Model doesn't have any output]**
-----------------------------Llama-2-13B---English-----------------------------
**Let's think step by step.**
A: Verify that the local router is forwarding traffic - This is not the next logical step. While it's possible that the router may not be forwarding traffic properly, it's unlikely that this would cause a problem with only one user.
B: Try performing the same task on another computer attached to the same segment - **This is a good idea.** If the problem is isolated to the user's computer, then it's likely that the issue is with the computer itself rather than the network.
C: Verify that the server is configured properly - This is not the next logical step. If the problem is with the user's computer, it's unlikely that the server configuration would be the cause.
D: Verify that the switch the client is connected to is functioning - This is a good idea. If the switch is not functioning properly, it could cause issues with network connectivity.
-----------------------------Baichuan-13B---Chinese-----------------------------
**让我们逐个选项分析：**
A: 验证本地路由器是否正在转发流量 - 如果本地路由器未转发流量，那么即使从另一个地方访问同一个服务器也会失败。然而，这个问题似乎不太可能是导致问题的唯一原因。
B: 尝试在连接到同一网段的另一台计算机上执行相同的任务 - **这个步骤对于诊断潜在的网络问题非常有用。**如果在另一台计算机上同样无法完成任务，说明问题不是因为特定的电脑而是整个网络环境导致的。
C: 验证服务器已正确配置 - 如果服务器未正确配置，可能会导致错误的信息，使其看起来好像无法从客户端访问。然而，这个问题不会影响整个网络。
D: 验证客户端连接的交换机是否正常工作 - 如果客户端连接的交换机出现故障或配置错误，可能会影响到客户端与服务器的通信。虽然这个选项看起来很有道理，但在很多情况下，这个问题并不一定是由于交换机故障导致的。
**因此答案是：** B

Figure 19: A failed English-answering example Generated by Baichuan-13B-Chat on few-shot evaluation under the CoT+SC setting in both English and Chinese.

# C   Annotation Guideline for OpsEval Categorization

## C.1   Overview

In the OpsEval project, we aim to categorize operational and maintenance tasks within the industry. This categorization process is pivotal for understanding the spectrum of tasks and the required abilities to address them effectively. The process involves two primary steps: automated screening using GPT-4 for initial topic modeling, followed by a manual review process involving domain experts.

## C.2   Task Categorization

### C.2.1   Objective

To categorize questions into one of eight distinct operational tasks based on industry relevance, task frequency, and significance within operational settings.

### C.2.2   Steps

1. **Review Initial Categorization**: Begin with the insights provided by GPT-4's topic modeling. Each question has been preliminarily categorized into one or more operational tasks.
2. **Understand Task Definitions**: Familiarize yourself with the details of the eight distinct tasks outlined in the provided Appendix. Each task has specific criteria and examples to guide your categorization.
3. **Assign Tasks**: For each question, decide which of the eight tasks it belongs to. A question should be categorized based on its core focus and the operational activity it pertains to.

4. **Justification**: Briefly justify your choice, especially if a question seems to fit into more than one category. Use the task definitions as a guide to support your decision.

### C.2.3 Detailed Task Categorizations

1. **General Knowledge**: Questions related to foundational concepts and practices in the Ops domain.
2. **Fault Analysis and Diagnostics**: Questions focusing on identifying and solving discrepancies or faults in systems or networks.
3. **Network Configuration**: Questions about optimal configurations for network devices to ensure efficient and secure operations.
4. **Software Deployment**: Questions dealing with the distribution and management of software applications.
5. **Monitoring and Alerts**: Questions on using monitoring tools to oversee system efficiency and setting up alert mechanisms.
6. **Performance Optimization**: Questions aimed at enhancing network and system performance.
7. **Automation Scripts**: Questions involving the creation of scripts to automate processes and reduce manual intervention.
8. **Miscellaneous**: Questions that do not fit into the above categories or involve elements from multiple categories.

### C.2.4 Task Categorization Template

> Question ID:
> Question: [Insert question text here]
> Assigned Task:
> Justification: [Provide a brief explanation for the task assignment here]

### C.2.5 Example for Task Categorization

> Question ID: 001
> Question: What steps should be taken to configure a firewall to prevent unauthorized access while allowing legitimate traffic?
> Assigned Task: Network Configuration
> Justification: This question specifically asks for optimal configuration strategies for a key network device (firewall) to ensure security and efficient operation, aligning perfectly with the 'Network Configuration' task.

## C.3 Ability Categorization

### C.3.1 Objective

To classify questions based on the required cognitive ability to answer them: Knowledge Recall, Analytical Thinking, or Practical Application.

### C.3.2 Steps

1. **Review Definitions**: Read the descriptions of the three abilities in the provided Appendix. Each ability category has distinct characteristics and examples.
2. **Evaluate Questions**: Assess the cognitive demand of each question. Consider what is primarily required to answer the question effectively: recalling information, analyzing data/situations, or applying knowledge in practical scenarios.
3. **Assign Ability Level**: Determine the most appropriate ability category for each question. Some questions may seem to require multiple abilities; choose the one that is most critical for addressing the core challenge of the question.

4. **Justification**: Provide a rationale for your categorization, especially for questions that may not clearly fit into a single category. Refer to the ability definitions to support your categorization.

### C.3.3 Detailed Ability Categorizations

1. **Knowledge Recall**: Requires recognizing and recalling core concepts and foundational knowledge.
2. **Analytical Thinking**: Demands deeper thought to dissect problems, correlate information, and derive conclusions.
3. **Practical Application**: Involves applying knowledge or analytical insights to provide actionable recommendations.

### C.3.4 Ability Categorization Template

Question ID:
Question: [Insert question text here]
Assigned Ability:
Justification: [Provide a brief explanation for the ability level assignment here]

### C.3.5 Example for Ability Categorization

Question ID: 002 Question: How would you optimize the performance of a network experiencing frequent bottlenecks?
Assigned Ability: Practical Application Justification: The question requires applying knowledge of network systems and performance optimization techniques to propose specific solutions, hence it falls under 'Practical Application'.

### C.4 General Guidelines

- **Consistency**: Strive for consistency in your categorization decisions. If similar questions are categorized differently, reassess your choices to ensure they align with the task and ability definitions.

- **Collaboration**: When in doubt, discuss challenging questions with fellow experts. Collaboration can help clarify ambiguities and refine the categorization process.

- **Documentation**: Keep detailed notes on your decisions, especially for questions that required significant deliberation. This documentation will be valuable for future reference and analysis.

By following these guidelines, you will contribute to a comprehensive and nuanced categorization of operational tasks and required abilities. This effort is crucial for enhancing our understanding of the operational landscape and the diverse skills professionals need to navigate it effectively.

## Checklist

1. For all authors...

    (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]

    (b) Did you describe the limitations of your work? [Yes] See Section 6.

    (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Section 6.

    (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

    (a) Did you state the full set of assumptions of all theoretical results? [N/A]

    (b) Did you include complete proofs of all theoretical results? [N/A]

3. If you ran experiments (e.g. for benchmarks)...

    (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]

    (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]

    (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] No, because the reported model performance is significant enough that p-values are not needed to verify the results.

    (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] In Appendix B.3, we present the compute resources used in evaluation.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

    (a) If your work uses existing assets, did you cite the creators? [Yes]

    (b) Did you mention the license of the assets? [Yes]

    (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]

    (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes]

    (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] We remove data which contains personally identifiable information or offensive content manually. See Section 3.2.

5. If you used crowdsourcing or conducted research with human subjects...

    (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [Yes] See Appendix C.

    (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

    (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [No] Since this is a community effort, all participants are voluntary and unpaid.