# OPSEVAL: A COMPREHENSIVE BENCHMARK SUITE FOR EVALUATING LARGE LANGUAGE MODELS' CA PABILITY IN IT OPERATIONS DOMAIN

#### Anonymous authors

Paper under double-blind review

#### 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. Existing NLP-related benchmarks (e.g., C-Eval, MMLU) can not guide the selection of suitable LLMs for Ops (OpsLLM), and current metrics (e.g., BLEU, ROUGE) can not adequately reflect the questionanswering (QA) effectiveness in the Ops domain. We propose a comprehensive benchmark suite, **OpsEval**, including an Ops-oriented evaluation dataset, an Ops evaluation benchmark, and a specially designed Ops QA evaluation method. Our dataset contains 7,334 multiple-choice questions and 1,736 QA questions. We have carefully selected and released 20% of the dataset written by domain experts in various sub-domains to assist current researchers in preliminary evaluations of OpsLLMs. We test over 24 latest LLMs under various settings such as selfconsistency, 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 0.9185 with human experts and is improved by 0.4471 and 1.366 compared to BLEU and ROUGE, respectively. Over the past one year, our dataset and leaderboard have been continuously updated.

034 035

037

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

#### 1 INTRODUCTION

039 The IT Operations (Ops) plays a crucial role in maintaining the efficient and stable operation of 040 information systems such as cloud computing, 5G networks<sup>1</sup> and financial information systems. As 041 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" 042 by Gartner (Lerner, 2017), this technique utilizes artificial intelligence to address tasks such as 043 anomaly detection, fault analysis, and performance optimization. In recent years, large language 044 models (LLMs) have witnessed significant advancements. The latest models, such as GPT-40 (Ope-045 nAI, 2024), GPT-4V (OpenAI, 2023b), Meta-Llama-3 (AI@Meta, 2024), and GLM-4 (Zeng et al., 046 2022), have demonstrated exceptional generalization and task-planning capabilities. As a result, 047 these models have provided numerous opportunities to enhance downstream domain-specific appli-048 cations. With its advanced text generation ability, LLM is well suited for Ops on tasks like question answering, information summarizing, and report analysis. Hereinafter, we refer to the LLM used 050 for Ops as **OpsLLM**, regardless of whether they have been optimized specifically for Ops.

 <sup>&</sup>lt;sup>1</sup>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.

While there are benchmarks for assessing general-purpose NLP-related capabilities, 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:

060 1) Sensitive data. The Ops data is primarily sensitive and proprietary to companies, with very few 061 publicly available data, making it difficult for any company to independently provide sufficient eval-062 uation data to ensure confidence in the test results. 2) Sub-domains. The Ops field spans many 063 sub-domains, like 5G communications, cloud computing, and bank transactions, each requiring a 064 mix of capabilities, or "tasks," such as network configuration or terminology explanation. The sheer number of sub-domains and tasks, combined with the absence of a systematic taxonomy, makes 065 classifying questions challenging. 3) Prompt sensitivity. Due to the relatively proprietary nature 066 of the Ops, existing LLMs have not undergone specialized supervised fine-tuning (SFT) for instruct 067 following within the Ops field, the evaluation results are more sensitive to prompt engineering. 068 Designing appropriate prompts for robust and accurate evaluation is challenging. 4) QA metric. 069 Existing metrics like BLEU focus on linguistic similarity between model output and reference an-070 swers, which often fails to capture true performance in Ops tasks. In Ops, it's essential to assess 071 whether the model's answers address key points in the reference and are supported by sufficient 072 evidence, reflecting the precise meanings of domain-specific terms. 073

- To address these issues, we propose **OpsEval**, a comprehensive benchmark suite for evaluating 074 LLMs' capability in the IT operations domain. First, to tackle the challenge of benchmark data 075 mostly being private, we initiated a community around AIOps, which has attracted dozens of com-076 panies to participate. We have selected 9 representative sub-domains from the community, allowing 077 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 079 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 081 field, we employ model-based pre-clustering and manual review to annotate eight tasks and three abilities. 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 learn-083 ing (ICL). Lastly, to address the inaccuracy of existing metrics in Ops QA evaluation, we design 084 FAE-Score, which evaluates model responses based on fluency, accuracy, and evidence, with each 085 criterion having its own dedicated assessment method. 086
- 087 The contributions of our paper are as follows: 1) We introduce **OpsEval**, the first bilingual multi-task 880 dataset in the Ops domain, covering 8 tasks and 3 abilities with 9,070 questions. To assist researchers in preliminary evaluating their OpsLLMs, we have carefully selected and released 20% of QAs 089 from our benchmark licensed under CC-BY-NC-4.0, with the remaining 80% of undisclosed data 090 preventing unfair evaluations due to data leakage (Wei & et.al., 2023) 2) Based on the dataset, we 091 introduce the OpsEval evaluation benchmark, conducting independent and robust evaluations with 092 various prompting techniques and a specifically designed evaluation metric, FAE-Score. Compared 093 to the commonly employed BLEU and ROUGE metrics, FAE-Score exhibits a more pronounced 094 congruence with the evaluations of human experts. Specifically, FAE-Score attains a correlation 095 coefficient 0.9175 with expert assessments, surpassing the coefficients of 0.6705 for BLEU and -096 0.3957 for ROUGE. 3) Based on the results of OpsEval evaluation, we provide key observations 097 and practical lessons to help domain practitioners make decisions such as whether existing models 098 are sufficiently applicable within a specific sub-domain, the necessity for fine-tuning and whether model quantization compromises the effectiveness. 099
- 100 101

102

#### 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 simi-



Figure 1: The framework of OpsEval

122 lar abilities rather than being confined to a particular domain. MMLU (Hendrycks et al., 2021) is a 123 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., 124 2022) employs seven distinct metrics in 42 unique scenarios, offering a comprehensive evaluation of 125 LLMs' capabilities across multiple dimensions. BIG-bench (Srivastava et al., 2022) comprises 204 126 tasks spanning a wide array of topics, with a particular focus on tasks deemed beyond the reach of 127 current LLMs. SEAL (AI, 2024b) features private, expert evaluations of leading frontiers models. 128 C-Eval (Huang et al., 2023) is a comprehensive Chinese evaluation suite designed to assess Chinese 129 LLMs' advanced knowledge and reasoning abilities rigorously. 130

**Domain-specific benchmarks** evaluate the abilities of LLMs to handle tasks in specific fields. 131 These benchmarks require LLMs to possess specialized knowledge in a specific domain and to 132 respond in a manner consistent with the cognitive patterns of that field. Despite the rapid progres-133 sion of LLMs in specialized domains, the evaluation metrics for these specific areas have received 134 less attention. FLUE (Shah et al., 2022) is an open-source comprehensive suite of benchmarks, 135 including new benchmarks across 5 NLP tasks in financial domain. MultiMedQA (Singhal et al., 136 2022) is an extensive medical question-answering dataset, with questions derived from professional 137 medical exams, research, and consultation records. CMB (Wang et al., 2023a) includes multi-choice 138 questions (CMB-Exam) and complex clinical questions based on real case studies (CMB-Clin). Ne-139 tOps (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 question-140 answering questions. However, they only focus on wired network operations and while the dataset is 141 released, they lack a benchmark that continuously updates the leaderboard. OWL (Guo et al., 2024) 142 introduces Owl-Instruct and Owl-Bench datasets for IT operations, along with methods like HMCE 143 for handling input length and a mixture-of-adapter for efficient tuning. However, it lacks a real-time 144 updated leaderboard and does not provide a well-designed evaluation for IT operations QA tasks. 145

146

121

### 3 OPSEVAL BENCHMARK

147 148

152

154

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.

153 3.1 DATA COLLECTION

Our benchmark questions have been collected from various sources; we summarize them into four categories: company materials, certification exams and Ops textbooks. 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 11 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.

162 Table 2: Overview of the question distribution in OpsEval by sub-domains, tasks and abilities. 163 (a) The number of questions in OpsEval, grouped (b) The distribution of different tasks and abilities 164 by their sub-domains. of questions in OpsEval.

Sub-domain	Source	Туре	Questions		Category	Percentage (%)
Wired Network	Operation Textbooks	MC	3901		Automation Scripts	3.3
5C Communication	Cartification Evans	MC	2615		Monitoring and Alerting	5.2
56 Communication	Certification Exams	QA	1162		Performance Optimization	5.3
Oracle Database	Company Materials	MC	497	Teals	Software Deployment	7.9
Log Analysis	Company Materials	QA	420	Task	Fault Analysis and Diagnostics	13.7
DevOps	Company Materials	QA	154		Network Configuration	29.0
Private Cloud	Company Materials	QA	150		General Ops Knowledge	20.2
Securities Info.	Company Materials	MC	91		Miscellaneous	15.5
Hybrid Cloud	Company Materials	MC	40		W 11 D 11	10.0
Financial IT	Company Materials	MC	40		Knowledge Recall	49.8
Total			9,070	Ability	Analytical Thinking Practical Application	$39.9 \\ 10.2$

175 Certification Exams. include knowledge assessments necessary for becoming an Ops staff and 176 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 177 sample questions from them as one of the sources for Ops questions. 178

179 **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 pro-181 vide experts with materials for question creation, and some books themselves also include a certain 182 number of exercises at the end of the chapters.

183

185

3.2 PREPROCESSING

186 We systematically carried out the preprocessing of our original data in the following stages: 187

**Deduplication:** Any repeated or highly similar questions are identified and removed to avoid redun-188 dancy in the test set. We calculate the cosine similarity of the question stems by bge-large-zh-v1.5 189 (Xiao et al., 2023) to detect duplicate questions and identify pairs of questions with a similarity 190 above a certain threshold (th=0.7). 191

192 **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 193 two parallel lists of empirical keywords in the question stems and the responses of GPT-3.5-turbo. 194 The keyword list can be found in the Appendix A.2. 195

196 **Question Categorization:** We devise a categorization that captures many tasks that professionals 197 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 20 199 tasks but had an imbalanced distribution. We then involved dozens of experts during the manual 200 review process to categorize the questions into eight tasks and three abilities. The distribution of the 201 questions across these eight tasks and three ability levels is shown in Table 2b, and the details of 202 each task and ability can be found in Appendix A.4. 203

Manual Review: In the manual review step, we asked Ops experts from the industry to inspect the 204 results of the previous three automated steps, including confirming duplicate and invalid questions 205 and examining the classification results of GPT-4. In our work, an expert is defined as an individ-206 ual with ten or more years of professional experience in their field, whether as an employee or a 207 researcher. Experts were also asked to drop the questions unrelated to the Ops field. We split the 208 dataset by n-folds and ensure each fold has at least two experts to review. As listed in Table 2a, this 209 quality enhancement process resulted in a refined test set of approximately 7,000 multi-choice and 210 2,000 question-answering questions. 211

212

214

#### 3.3 EVALUATION SETTINGS 213

**Multi-choice questions** offer a structured approach with definitive answers. These questions are 215 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. We evaluate question-answering tasks using a metric designed
 specifically for OpsEval, called FAE-Score, which is explained in detail in the subsequent section.
 Additionally, we perform expert evaluations and calculate BLEU (Papineni et al., 2002), ROUGE
 (Lin, 2004) and RAGAS (Es et al., 2024) scores for comparison purposes, as reference to validate
 the accuracy of FAE-Score.

- We use the same three criteria to evaluate the responses of various models for both FAE-Score and Expert Evaluation:
- 226
- 227 228

229

230

231

• **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.

In 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.

Prompting Techniques. We use various settings to evaluate LLMs on OpsEval to get a comprehen sive 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 24 LLMs can be found in Table 6 in Appendix C.1.

246 247

3.4 FAE-Score

Figure 3 shows the basic pipeline of our designed
QA metric, FAE-Score. Here, we elaborate each
evaluation methodology of each criterion.

Fluency. In Ops settings, the fluency of a model's output is crucial because the results are intended for 253 human consumption by technical personnel. Unlike 254 other generic benchmarks, the tasks in the Ops domain often require clear and unambiguous commu-256 nication, as the model's outputs may guide decision-257 making in real-world scenarios. Therefore, ensuring 258 high fluency in responses is not just a matter of lan-259 guage quality but a critical factor for task comple-260 tion and user comprehension. To evaluate fluency 261 in model outputs, we adapted the scoring rubrics methodology mentioned in Kim et al. (2024). We 262



Figure 2: Scoring rubrics for Fluency metric.

use Qwen2-72B-Instruct as the evaluation model, for its strong performance in general language
 generation (QwenLM, 2023) and its consistent multilingual capabilities without significant degra dation. We assess the fluency of various model outputs, scoring them based on grammar, coherence,
 clarity, appropriateness of style, and answer completeness, as shown in the Figure 2.

Accuracy. Traditional metrics such as BLEU and ROUGE fall short in this vertical domain because
 they often fail to capture the key factual content within long-form responses. This results in inflated
 scores due to irrelevant word matches, making these metrics insufficient for accuracy evaluation in
 the highly specialized and knowledge-driven Ops context. To address these shortcomings, we take



Figure 3: The FAE-Score pipeline.

inspiration from Es et al. (2024), using a keyword extraction method to evaluate the accuracy of model outputs. A judge model (OpenAI, 2023a) is then employed to match the keywords from the model's response with the keywords from the standard answer. The final accuracy score is calculated by determining the F1-Score, which balances precision and recall for the matched keywords.

Accuracy = 
$$2 \cdot \frac{P \cdot R}{P + R}$$
,  $P = \frac{\#\text{Matched Keywords}}{\#\text{Keywords in Model Output}}$ ,  $R = \frac{\#\text{Matched Keywords}}{\#\text{Keywords in Ground Truth}}$  (1)

**Evidence.** Model responses must not only be accurate but also well-supported by relevant, authoritative information. To evaluate the evidence behind a model's response, we implement a ROUGEbased method to measure the overlap between the generated output and the content of related documents retrieved through similarity search. We used bge-large-zh (Xiao et al., 2023) for document embedding and FAISS (Douze et al., 2024) for similarity search. By retrieving documents that closely match the question, we can assess whether the model's response appropriately references or aligns with this external information. We use ROUGE, as a recall-oriented metric, captures how much of the content in the relevant documents is reflected in the model's output. This ensures that the model does not simply generate plausible-sounding answers but grounds its responses in factual evidence from trusted sources.

$$Evidence = ROUGE_{Recall}(R, D) = \frac{\#Overlapping Words}{\#Words in D}$$
(2)

3.5 OPEN-SOURCE POLICY

We have released 20% of the OpsEval dataset to the public to foster contributions from the Ops 305 community and support research. To ensure balanced distribution, this subset was randomly sam-306 pled from each data source and sub-domain in proportion to their respective weights. Additionally, 307 for questions involving proprietary company data, we carefully reviewed and modified the content 308 to remove any sensitive information. This sample dataset provides researchers with insights into the 309 types and topics of questions expected in the benchmarks, allowing them to better understand the 310 scope of the evaluation. The sampled dataset also enables model developers to conduct local eval-311 uations of their models, facilitating faster iterations. Moreover, this dataset can serve as a seed for 312 generating QA pairs through automatic QA generation algorithms (Wang et al., 2023c), contribut-313 ing to the growth of Ops-specific data for future model development. While this subset is available for users' self-evaluation, the complete dataset remains undisclosed. By ensuring that the test set 314 answers are not leaked, we guarantee the reliability and non-leakage of the OpsEval benchmark. 315

316 317

318

320

281

282

283

284

287

289

290

291

292

293

295

296

297

298 299 300

301 302 303

304

- 4 RESULT ANALYSIS
- 319 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 4. Results of the other sub-domains and settings are shown in Appendix C.4.
 While closed source models like GPT-4 and Claude-3-Opus performs well on the OpsEval benchmark, open-sourced LLMs yield generally worse evaluation results than those in general domains

336

337

338

339

340 341 342

343

344 345

347

348

349

350

363

364

366



Figure 4: 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.



Figure 5: 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.

 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. Recent open-sourced models like
 Qwen2-72B-Chat, exhibit competitive performance in multi-choice questions, thanks to their finetuning process and the quality of their training data. 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
- 367 To investigate how LLMs perform in each Ops sub-domain and each task, and to what extent they 368 possess the general abilities, we summarize the result of different parameter-size groups of LLM 369 and plot them on three radar charts in Figure 5. Regarding the eight tasks we tested, LLMs yield 370 higher accuracy in General Knowledge tasks, while their performance drops and varies drastically in 371 highly specialized tasks like Automation Scripts and Network Configuration, reflecting the impact 372 of specialized corpus and domain knowledge on the performance of LLMs. By grouping LLMs 373 by their parameter size, we find that while LLMs with 10B-30B parameters have higher accuracy 374 in their best cases compared with LLMs with no more than 10B parameters, different 10B-20B 375 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 376 LLMs across the 8 tasks and 3 abilities in Figure 6. By examining the focus areas across different 377 categories, we have identified key research targets for capability training.

Software Deployment Fault Analysis Network and Diagnostics Configuration General Ops Knowledge Automation Scripts Monitoring and Alerting Performance Optimization Miscellaneous 0.625 Knowledge Recall 0.56 0.600 0.575 Analytical Thinking 0.550 0.64 0.54 0.54 0.56 0.54 0.56 0.525 0.500 Practical 0.57 Application 0.475

Figure 6: Heatmap of failure case distribution regarding tasks and abilities. The values represent the proportion of failure cases across all LLMs; yellower areas indicate higher failure rates.

Table 3: LLMs' performance on English network operations question-answering problems.

Model	ROUGE(%)	BLEU(%)	RAGAS(0-10)	Fluency		Accuracy		Evidence		FAE-Total	
	100001(10)	DELC(,n)		FAE	Expert	FAE	Expert	FAE	Expert	FAE	Expert
GPT-3.5-turbo	12.26	6.78	9.23	9.38	9.12	8.06	9.65	6.21	8.11	23.65	26.88
LLaMA-2-70B	7.74	4.2	6.04	8.69	8.25	7.71	8.79	9.08	8.98	25.48	26.02
LLaMA-2-13B	4.98	3.43	8.23	8.47	9.84	7.32	9.34	8.81	7.27	24.60	26.44
Chinese-Alpaca-2-13B	3.25	1.85	5.32	5.53	8.05	6.99	7.95	6.23	6.23	18.75	22.24
Baichuan-13B-Chat	4.76	0.35	7.93	7.16	7.98	8.71	7.84	6.66	7.31	22.53	23.13
Qwen-7B-Chat	11.82	4.33	4.92	7.63	5.82	6.42	7.27	6.57	5.37	20.62	18.47
ChatGLM2-6B	9.71	5.07	5.32	5.12	7.96	6.41	6.39	6.14	4.32	17.67	18.67
InternLM-7B	13.27	0.54	6.21	4.99	5.16	5.00	4.90	4.75	4.28	14.74	15.77
Chinese-LLaMA-2-13B	9.19	0.24	7.34	6.98	4.64	5.29	6.32	4.63	8.34	16.90	17.88

**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.

409 410

378

379 380

381

382

384

385

386

387

388

389

396 397

399 400

401

402

403

404

405

406

407

408

411

412 413

#### 4.3 PERFORMANCE ON QUESTION-ANSWERING

414 415

Table 3 presents the evaluation results of 200 question-answering English questions across four met-416 rics: ROUGE, BLEU, RAGAS, FAE-Score, and Expert-Evaluation. To gain more insight into how 417 different metrics perform in QA evaluation, we use Figure 20 (see in Appendix C.9.2) as a case 418 analysis. While BLEU and ROUGE are efficient in natural language comparison, they lack seman-419 tic information to determine which part of the context is more important than others. Knowing that 420 a given benchmark evaluates QA based on BLEU/ROUGE, there is an obvious way to trick the met-421 ric: repeat patterns occurring in the question, gaining a higher possibility to match some patterns 422 in the reference answer. Due to their lack of semantic information related to Ops and the potential 423 hack, traditional metrics like BLEU are unsuitable for specialized benchmarks. Instead, with specialized prompting and seperately designed methodology for each criterion (Fluency, Accuracy and 424 Evidence), FAE-Score can comprehensively evaluate models' QA performance, with the Accuracy 425 metric picking up those important keywords and not be influenced by repeated words that contain 426 no useful information, and the Evidence metric checking the recall of relevant supporting contents. 427 In Section 5, we discuss the alignment between different metrics and expert evaluation, validating 428 the effectiveness of FAE-Score in automated QA evaluation within the Ops domain. 429

Insight: In specialized domains, Ops specifically, traditional NLP metrics like BLEU and ROUGE
 cannot comprehend the key components in the reference answer, resulting their evaluation lacking practical significance. FAE-Score is suitable for large-scale qualitative evaluations in the Ops field.

433 434 425	(a) Measurement of potential test data leak- age during the training of LLM.						) Pearson valuation n	correlation netrics and	n coefficie Automate	ents betweed metrics.	een Expert- Total is the
430	Dataset	$L_{test}$	$L_{ref}$	$\Delta L$	$\geq 0?$	-		icy, Accur	acy, and E	vidence.	
497	Alpaca	1.9940	2.3542	-0.3602	x		Metric	Total	Flu.	Acc.	Evi.
437	Alpaca-GPT4	1.4988	1.7636	-0.3910	×		ROUGE	-0.44734	-0.49207	-0.40889	-0.31821
438	CEval	2 5708	2 3099	0 2608	1		BLEU	0.47139	0.46369	0.55330	0.05977
439	MMLU	2.5475	2.1898	0.3577	1		RAGAS	0.57169	0.40029	0.51151	0.41928
440	OpsEval	2.9854	2.6280	0.3050	1	_	FAE-Score	0.91848	0.54757	0.81523	0.58160

Table 4: Validation results.

#### 4.4 PERFORMANCE ON DIFFERENT QUANTIZATION PARAMETERS

We conducted experiments on different quantized versions of LLaMA-2-70B and obtained various results and conclusions. For detailed results, please see Appendix C.5. 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.

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

449 450

432

441

442 443

444

445

446

447 448

451 452

#### 5 VALIDATION

- 453 454
- 455

#### 5.1 BENCHMARK LEAKAGE TEST

456 For the fairness of a benchmark suited for LLM, avoiding potential bias emerging from test set leak-457 age is necessary. We adapted the methodology from Wei & et.al. (2023) to perform a leakage test on 458 OpsEval's dataset. We evaluate the LLM loss on samples from different datasets for several LLMs 459 and calculate the average loss. For each dataset, we compare LLM loss on the test split  $(L_{test})$  and 460 a specially curated reference set ( $L_{ref}$ ) generated by GPT-4, designed to mimic the testing dataset. 461 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 462 and accuracy.<sup>2</sup> We define a key metric:  $\Delta L = L_{test} - L_{ref}$ , with a threshold of  $\Delta L < 0$  indi-463 cating potential test data leakage. A negative  $\Delta L$  suggests that the LLM's lower  $L_{test}$  comes from 464 overfitting the test set rather than understanding the questions, indicating potential leakage. Table 4a 465 shows the results of leakage measurement. In addition to the two standard evaluation benchmarks 466 (CEval (Huang et al., 2023) and MMLU (Hendrycks et al., 2021)), we conducted the same exper-467 iments on the alpaca dataset (Taori et al., 2023) and the Alpaca-GPT4 dataset (Peng et al., 2023), 468 which is likely used in the pre-training of large models, using its  $\Delta L$  as reference. This demonstrates 469 the unbiased nature and non-leakage of the OpsEval test set. The models used in the leakage test are 470 listed in Appendix C.1.

471 472

473 474

#### 5.2 EXPERT ALIGNMENT OF FAE-SCORE

Table 4b shows the correlation coefficients between various automated scoring metrics (ROUGE, 475 BLEU, RAGAS, and FAE-Score) and Expert-Evaluation criteria. The results indicate that ROUGE 476 and BLEU scores often misalign with Expert-Evaluation. This misalignment occurs because LLMs 477 with poor performance may generate keywords that boost ROUGE and BLEU scores, while stronger 478 LLMs might receive lower scores due to different wording from standard answers. While RAGAS 479 (Es et al., 2024) aligns better with experts than ROUGE and BLEU, there is still a gap between 480 its scoring rankings for different models and expert judgement standards. In contrast, FAE-Score 481 rankings closely match Expert-Evaluation, particularly with the Accuracy metric. This suggests that 482 FAE-Score is more reliable in assessing the factual accuracy of LLMs' outputs. Notably, GPT-4's 483 performance in factual accuracy is reflected in its strong alignment with the Accuracy metric.

<sup>&</sup>lt;sup>2</sup>For a case example, please see Appendix C.8

## 486 6 DISCUSSION

## 488 6.1 AUTOMATED QA GENERATION

During the data collection process, we explored automating question-answer generation. Initially,
we sampled QA pairs and manually evaluated their accuracy and domain relevance. Later, we
utilized representative examples for few-shot learning, enabling GPT to generate and evaluate QA
pairs automatically based on predefined criteria.

494 Recognizing that most existing benchmarks focus primarily on simple knowledge-based questions, 495 we designed various task-specific templates to address this limitation. These templates require the model to complete specific fields within the template using the provided knowledge content, rather 496 than generating entire questions and answers. This prompt engineering approach allows us to gen-497 erate detailed and context-specific Ops tasks based on extensive operational knowledge while im-498 proving the model's instruction-following ability. By focusing on field-level completion, the overall 499 structure of the QA remains consistent and accurate. In the appendix, we provide the prompt tem-500 plate used for automatic QA generation (Figure 11), along with some task cases illustrating their 501 application (Figure 12). This approach ensures a more diverse and comprehensive evaluation of 502 model capabilities while maintaining the relevance and quality of generated tasks.

503 504

505

6.2 FUTURE WORK

506 **Comprehensive Error Analysis.** To better understand the limitations of large models in Ops ques-507 tion answering, we will further look into the failure cases and identify common error modes, includ-508 ing lack of domain knowledge, hallucinations, inaccurate reasoning, and overconfidence in incorrect 509 answers. We believe this detailed error analysis provides a clearer picture of the challenges faced by models and informs future research directions to address these issues. Dataset Scale and Real-510 World Data. While privacy constraints limit real-world company data, our ongoing collaborations 511 aim to expand the dataset with practical scenarios. Expanding the dataset with real-world scenarios 512 remains a key focus, while the benchmark prioritizes robust evaluation over dataset scale. Agent 513 and RAG Introduction: The inclusion of agents and Retrieval-Augmented Generation (RAG) tech-514 niques is constrained by the current large models' lack of foundational knowledge in operations. Our 515 leaderboard will incorporate more complex tasks once open-source models possess sufficient oper-516 ational capabilities.

517 518

519

#### 7 CONCLUSION

In this paper, we introduced **OpsEval**, the first comprehensive Ops benchmark suite designed for 520 evaluating the performance of large language models (LLMs) in IT operations. We established a ro-521 bust evaluation framework encompassing a wide range of sub-domains and tasks within Ops through 522 rigorous data collection from multiple sources and meticulous preprocessing steps. Our benchmark 523 includes a carefully selected set of 9,070 questions, which we have partially released to aid initial 524 evaluations while protecting the integrity of the remaining dataset. It has undergone experiments 525 in data leakage detection, ensuring its reliability. Our observations, supported by quantitative and 526 qualitative results, highlight the need for a balanced approach to selecting fundamental models, con-527 sidering both performance and robustness. During the QA evaluation, the FAE-Score emerges as a 528 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 abili-529 ties provides researchers with crucial insights and prospects for future breakthroughs. The identified 530 flexibility within the OpsEval framework presents opportunities for future exploration. This bench-531 mark's adaptability facilitates the seamless integration of additional fine-grained tasks, providing a 532 foundation for continued research and optimization of LLMs tailored for Ops. 533

- 534
- 535 **REFERENCES**

 536
 537
 538
 CodeFuse AI. Codefuse-devops-model, 2024a. URL https://github.com/codefuse-ai/ CodeFuse-DevOps-Model. Accessed: 2024-06-06.

539 Scale AI. Seal leaderboards, 2024b. URL https://scale.com/leaderboard. Accessed: 2024-06-03.

540 AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/ 541 blob/main/MODEL\_CARD.md. 542 Baichuan. Baichuan 2: Open large-scale language models. arXiv preprint arXiv:2309.10305, 2023. 543 URL https://arxiv.org/abs/2309.10305. 544 Baidu. baidu/ernie-bot-4.0, September 2024. URL https://cloud.baidu.com/doc/ 546 WENXINWORKSHOP/s/clntwmv7t. 547 548 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, 549 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 550 2021. 551 552 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-553 Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024. 554 555 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: 556 General language model pretraining with autoregressive blank infilling. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 320-335, 2022. 558 559 Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. RAGAs: Automated evalu-560 ation of retrieval augmented generation. In Nikolaos Aletras and Orphee De Clercq (eds.), Pro-561 ceedings of the 18th Conference of the European Chapter of the Association for Computational 562 Linguistics: System Demonstrations, pp. 150–158, St. Julians, Malta, March 2024. Association 563 for Computational Linguistics. URL https://aclanthology.org/2024.eacl-demo. 564 16. 565 Gemma\_Team, Thomas Mesnard, and et.al. Gemma: Open models based on gemini research and 566 technology, 2024. 567 568 Hongcheng Guo, Jian Yang, Jiaheng Liu, Liqun Yang, Linzheng Chai, Jiaqi Bai, Junran Peng, Xi-569 aorong Hu, Chao Chen, Dongfeng Zhang, xu Shi, Tieqiao Zheng, liangfan zheng, Bo Zhang, 570 Ke Xu, and Zhoujun Li. OWL: A large language model for IT operations. In The Twelfth Interna-571 tional Conference on Learning Representations, 2024. URL https://openreview.net/ forum?id=SZOQ9RKYJu. 572 573 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 574 Steinhardt. Measuring massive multitask language understanding. Proceedings of the Interna-575 tional Conference on Learning Representations (ICLR), 2021. 576 577 Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, 578 Chuancheng Lv, Yikai Zhang, Jiayi Lei, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. arXiv e-prints, pp. arXiv-2305, 2023. 579 580 InternLM\_Team. InternIm: A multilingual language model with progressively enhanced capabilities. 581 https://github.com/InternLM/InternLM, 2023. 582 583 Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham 584 Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Prometheus 2: An open source language 585 model specialized in evaluating other language models, 2024. 586 Andrew Lerner. Aiops platforms-gartner, 2017. 587 588 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian 589 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language 590 models. arXiv e-prints, pp. arXiv-2211, 2022. 591 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In Text Summarization 592 Branches Out, pp. 74-81, Barcelona, Spain, July 2004. Association for Computational Linguis-593

tics. URL https://aclanthology.org/W04-1013.

612

624

632

633

634

635

- Yukai Miao, Yu Bai, Li Chen, Dan Li, Haifeng Sun, Xizheng Wang, Ziqiu Luo, Dapeng Sun, Xiuting Xu, Qi Zhang, Chao Xiang, and Xinchi Li. An empirical study of netops capability of pre-trained large language models. *CoRR*, abs/2309.05557, 2023. URL https://doi.org/10.48550/arXiv.2309.05557.
- 598 599 OpenAI. Chatgpt: Optimizing language models for dialogue. OpenAI Blog, 2022. URL https: //openai.com/blog/chatgpt/.
- <sup>601</sup> OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023a.
- OpenAI. Gpt-4v(ision) system card, 2023b. URL https://cdn.openai.com/papers/ GPTV\_System\_Card.pdf.
- 605 OpenAI. Hello gpt-40, 2024. URL https://openai.com/index/hello-gpt-40/. Accessed: 2024-06-01.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL https://aclanthology.org/P02-1040.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. arXiv preprint arXiv:2304.03277, 2023.
- 615 QwenLM. Qwenlm/qwen-7b, September 2023. URL https://github.com/QwenLM/
   616 Qwen-7B.
- Raj Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. When FLUE meets FLANG: Benchmarks and large pretrained language model for financial domain. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2322–2335, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.148. URL https://aclanthology.org/2022.emnlp-main.148.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv e-prints*, pp. arXiv–2206, 2022.
  - Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford\_alpaca, 2023.
- Hugo Touvron and et.al. Llama 2: Open foundation and fine-tuned chat models, 2023.
- Kidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao,
   Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, et al. Cmb: A comprehensive medical
   benchmark in chinese. *arXiv e-prints*, pp. arXiv–2308, 2023a.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023b.
- Yizhong Wang, Yeganeh Kordi, and Swaroop et al. Mishra. Self-instruct: Aligning language models
  with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–13508, Toronto, Canada,
  2023c. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.754. URL
  https://aclanthology.org/2023.acl-long.754.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
  - Tianwen Wei and et.al. Skywork: A more open bilingual foundation model, 2023.
  - Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. C-pack: Packaged resources to advance general chinese embedding, 2023.
  - Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
  - A DETAILS OF OPSEVAL BENCHMARK
  - A.1 INFORMATION ON THE COMPANIES AND EXPERTS PARTICIPATING 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/
Inspur	Ops service/tool provider	https://www.inspur.com/
Total	11	

 Table 5: Information of companies collaborating in OpsEval

Table 5 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 PROMPT FOR GPT-4 CATEGORIZATION

Figure 7 shows the prompt for GPT-4 initial categorization.

A.4 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.

I need your help in analyzing a multi-choice question, de	etermine the domain and the
task type it belongs to.	
Domains: When classifying the domain, be specific, div	e deeper into domains such as:
Database/Network Operations	
Task Types: For the task type, consider categories like: ]	Monitoring and Alerts,
Performance Optimization	C ·
Summary your response as JSON format: {"domain"	': "specific domain", "task":
"specific_task_type"}	, in the second s
Figure 7: The prompt for GPT-4 initial	categorization
Note I C for the prompt for OT 1-4 million	
• <i>Network Configuration</i> revolves around suggesting optimal like routers, switches, and firewalls to ensure their efficient	at and secure operations.
• <i>Software Deployment</i> deals with the dissemination and n throughout the network or system, verifying their correct	nanagement of software applications installation.
• <i>Monitoring and Alerts</i> harnesses monitoring tools to super and implements alert mechanisms to notify administrators	ervise network and system efficiency sof emerging issues.
• <i>Performance Optimization</i> is centered on refining the netw and recognizing potential enhancement areas.	ork and system for peak performance
• Automation Scripts involves the formulation of automatic	on scripts to facilitate processes and
decrease manual intervention for administrators.	1 1
• <i>Miscellaneous</i> comprises tasks that do not strictly adhere	to the aforementioned classification
or involve a combination of various tasks.	
A.5 ABILITY LEVELS OF QUESTIONS	
Different questions require different levels of ability to answer. nto 3 categories. The details of each ability are as follows:	We classify all questions in OpsEva
1. <i>Knowledge Recall:</i> Questions under this category primaril	y test a model's capacity to recognize
and recall core concepts and foundational knowledge. S	such questions are akin to situation
where a professional might need to identify a standard p	rocedure or recognize a well-know
issue based solely on previous knowledge.	11 (77)
2. Analytical thinking: These questions demand more than m	here recall. They necessitate a deepe
mation and derive a coherent conclusion. It mirrors real-	world scenarios where professional
troubleshoot complex issues by connecting various dots	and leveraging their comprehensiv
understanding.	and to for going men comprehensiv
3. Practical Application: These questions challenge a mod	el's ability to apply its foundationa
knowledge or analytical conclusions to provide actionabl	e recommendations for specific sce
narios. It epitomizes situations where professionals are ex	spected to make decisions or suggest
solutions based on in-depth analysis and expertise.	
igure 8 illustrates examples in our question set, shedding ligh	t on our classification methodology.
A.6 PROMPT AND FORMATTING OF QUESTIONS	
Zigure 9 illustrates examples of the questions often our propriet	assing ningling
rigure 8 mustrates examples of the questions after our preproc	essing pipeline.
A.7 AN EXAMPLE OF SUBJECTIVE QUESTIONS	
A saved subjective question in OnsEval is presented in Figure	9 which contains not only the roy
nuestion but also its type of task.	, which contains not only the fav
As snown in Figure 10, we combine the task and ability of eac	n question with the question itself a

750	
757	Which of the following represents quantifying data moved from one host to another within a
758	specific time frame?
759	A: Reliability B: Response time C: Throughput D: Jitter
760	Answer: C
761	amount of time
762	Task: Performance Optimization Ability: Knowledge Recall
763	
764	details from a DHCPv6 server?
765	A: ipv6 nd ra suppress B: ipv6 dhop relay destination
766	Answer: D
767	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
768	Task: Automation Scripts
769	Ability. Analytical miniking
770	Question: You receive a call from a user experiencing difficulties connecting to a new VPN. What is the initial step you should take?
771	A: Find out what has changed. B: Reboot the workstation.
772	Answer: D
773	Analysis: Since this is a new connection, you need to start by troubleshooting and identify the symptoms and notential causes
774	Task: Fault Analysis and Diagnostics
775	Ability: Practical Application
776	Figure 8: Three examples of the processed questions
777	
778	
779	Question: You have a router interface with an IP address of 102 168 102 10/20. What is the broadcast
780	address that the hosts on this LAN will utilize?
781	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?
782	Keypoint: 192.168.192.15 然安亚古: 102.169.102.15
783	$= 2.5 \times 10^{-10}$
784	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.
785	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。
786	Task: Network Configuration
787	任务:网络配置
788	Ability: Analytical Thinking
789	(能力:推理)
790	
791	Figure 9: An example of the saved subjective questions
	Figure 9: An example of the saved subjective questions
792	Figure 9: An example of the saved subjective questions
792 793	Figure 9: An example of the saved subjective questions
792 793 794	Figure 9: An example of the saved subjective questions A subjective question in OpsEval Overtical You have a router interface with an IB address of 102 100 100 10/20. When is the how the set
792 793 794 795	Figure 9: An example of the saved subjective questions 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?
792 793 794 795 796	Figure 9: An example of the saved subjective questions 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?。
792 793 794 795 796 797	Figure 9: An example of the saved subjective questions 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 在本、网络配要
792 793 794 795 796 797 798	Figure 9: An example of the saved subjective questions 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
792 793 794 795 796 797 798 799	Figure 9: An example of the saved subjective questions 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
792 793 794 795 796 797 798 799 800	Figure 9: An example of the saved subjective questions 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 Promet
792 793 794 795 796 797 798 799 800 801	Figure 9: An example of the saved subjective questions 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
792 793 794 795 796 797 798 799 800 801 801	Figure 9: An example of the saved subjective questions 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
<ul> <li>792</li> <li>793</li> <li>794</li> <li>795</li> <li>796</li> <li>797</li> <li>798</li> <li>799</li> <li>800</li> <li>801</li> <li>802</li> <li>803</li> </ul>	Figure 9: An example of the saved subjective questions 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?
792 793 794 795 796 797 798 799 800 801 802 803 804	Figure 9: An example of the saved subjective questions         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?
792         793         794         795         796         797         798         799         800         801         802         803         804         805	Figure 9: An example of the saved subjective questions 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? IOMENTIAL Configuration Generation Generation Description Descri
792         793         794         795         796         797         798         799         800         801         802         803         804         805         806	Figure 9: An example of the saved subjective questions 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? @EXTENDED Configuration Gamma Analytical Thinking ####################################
792         793         794         795         796         797         798         799         800         801         802         803         804         805         806         807	Figure 9: An example of the saved subjective questions 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?  Maswer 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?  Maswer 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?  Maswer 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?  Maswer 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?  Maswer the Reasoning the provide the provid
792         793         794         795         796         797         798         799         800         801         802         803         804         805         806         807         808	Figure 9: An example of the saved subjective questions         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?         jdm: BtaBmachart and the hosts on this LAN will utilize?         jdm: BtaBmachart and the host on this LAN will utilize?         jdm: BtaBmachart and the host on this LAN will utilize?         jdm: BtaBmachart and the host on this LAN will utilize?         jdm: BtaBmachart and the host on this LAN will utilize?         jdm: BtaBmachart and the host on the law the host on the law will utilize?         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 host on this LAN will utilize?         jdmAshmachart and the address of 192.108.192.10/29. What is the broadcast address that the hosts on this LAN will utilize?         jdmAshmachart and host intervention.         You have a router interface with an IP address of 192.108.192.10/29. What is the broadcast address that the host on this LAN will utilize?         jdmAshmachart and the address of 192.108.192.10/29. The host on the ho

839

840

841 842

843

844 845

846 847

848

0.0	
811	
812	You are an operations expert, and your task is to generate a question that adheres to the
813	given template or follows a similar format, based on the provided knowledge content.
814	The question template is as follows:
815	
816	{question_template}
817	
818	Here, {{}} represents parameters that need to be filled.
819	The knowledge points for generating the question are as follows:
820	
821	{context}
822	
823	when generating the question, you need to construct a business or operations scenario based on the knowledge content, and then use that scenario to populate the required
824	fields.
825	
826	When creating the question, you need to adhere to the following constraints:
827	{constraints}
828	Poturn a ISON object containing the required fields from the template, as well as an
829	answer field and an explanation field. The answer field should contain the answer to the
830	question, and the explanation field should provide reasoning for the answer, explaining
831	why it is correct. The terms used in your question and answer must match the given
832	knowledge content, and you should not invent new terminology.
833	The reference format is as follows:
834	
835	{json_example}
836	
837	Your returned content should not start with ```json. Return the JSON object directly.
000	

#### Figure 11: Prompt template for automated QA generation

### **B** AUTOMATED QA GENERATION

Figure 11 shows the prompt template we used for automated QA generation experiment. Figure 12 shows some automatically generated QAs, their task description, template and example question.

C ADDITIONAL DETAILS OF EXPERIMENTS

#### C.1 DETAILED INFORMATION OF LLMs EVALUATED

849 GPT-4 (OpenAI, 2023a) is a large multimodal model (accepting image and text inputs, emitting text 850 outputs) that, while less capable than humans in many real-world scenarios, exhibits human-level 851 performance on various professional and academic benchmarks. It is recognized as the strongest 852 lanuage model currently. ChatGPT (OpenAI, 2022) is an earlier AI-powered language model devel-853 oped by OpenAI which is built upon GPT-3.5. We use the GPT-3.5-turbo version in our experiments. 854 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 855 Chinese. We evaluate three weights (70B, 13B and 7B as shown in 6) of LLaMA 2. 856

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 (Baidu, 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 (QwenLM, 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 6) of Qwen-Chat. Baichuan2-13B-Chat (Baichuan, 2023) is aligned chat model based

865		
866		
867	Business Information Reasoning	Tool Selection
868 869	Task definition: The model needs to analyze the provided business information and solve user problems based on the given business knowledge and reasoning rules.	Task definition: The model should select and use the appropriate tool to solve user issues based on the provided operations scenario, tool descriptions, and user queries.
870 871 872 873 874	Template: You are a core network operations engineer. Your task is to strictly follow the reasoning rules to provide an analysis result. Keep the reasoning process as concise as possible. Input Information: (input_info) Business Knowledge: (business_knowledge)	Template: You can use various user-defined tools to solve the given user issues. Your task is to resolve the user's question based on the tool description. Tool Description: [tool_description] User Question: [user_question]
876 877	You need to answer in the specified format, filling in the content according to the reasoning rules. The response format should be in JSON and include two fields: reasoning process and analysis result.	Example: Tool Description: ( "name": "DSP_PAENODE", "
878 879 880 881 882 883	Example: Input Information: Work Order Content: [HSS Subscription Conclusion] Voice service unavailable: Incoming call lock, unable to receive voice calls. Data service unavailable: Locked GPRS, unable to use 2G data service; Locked EPS, unable to use 4G data service. Conclusion: Analyze HSS	line of the second sec
884 885 886 887	Business Knowledge: Focus only on whether there are any issues with data services in the work order information; ignore other aspects. If the HSS subscription conclusion indicates that the data service is normal, output "Analyze PGW." If it indicates the data service is abnormal, output "Analyze HSS." Based on the and, analyze the work order content in the and output strictly according to the reasoning rules. The output should be:	"arguments": {     "meid": "String Network Element ID",     "nodeName": "String NODE Name"     },     "results": "NODE-related information"     }     User Question:     Retrieve the host name of the NODE with Network Element ID 232 and NODE name 2021-
888		8-88-49-6fb9, and provide the parameters corresponding to the query in JSON format.

Figure 12: Some automatically generated QAs, their task description, template and example question

Table 6: 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.

899	Model	Creator	#Parameters	Access	License
900	GPT-4/GPT-3.5-turbo	OpenAI	undisclosed	API	Proprietary
902	ERNIE-Bot-4.0	Baidu	undisclosed	API	Proprietary
903	GLM4/GLM3-turbo	Tsinghua Zhipu	undisclosed	API	Proprietary
904	Meta-LLaMA-3	Meta	8B	Weights	Llama 3 Community
905	LLaMA-2	Meta	7/13/70B	Weights	Llama 2 Community
906	Qwen-Chat	Alibaba Cloud	7/14/72B	Weights	Qianwen LICENSE
907	Qwen1.5-Chat	Alibaba Cloud	14B	Weights	Qianwen LICENSE
908	InternLM2-Chat	Shanghai AI Laboratory	7/20B	Weights	Apache-2.0
909	DevOps-Model-Chat	CodeFuse	14B	Weights	Apache-2.0
910	Baichuan2-Chat	Baichuan Intelligence	13B	Weights	Apache-2.0
911	ChatGLM3	Tsinghua Zhipu	6B	Weights	Apache-2.0
912	Mistral	Mistral	7B	Weights	Apache-2.0
913	Gemma	Google	2/7B	Weights	Gemma license
914	Claude-3-Opus	Anthropic	undisclosed	API	Proprietary
915	Qwen2-Instruct	Alibaba Cloud	7/72B	Weights	Qianwen LICENSE
916					

918	Tab	Table 7: GPTQ models for LLaMA-2-70B												
919	Model	Size	#CPTO Dataset	Dise										
920	Widder	Size		Disc										
921	LLaMA-2-70B	140GB	/	Raw LLaMA-2-70B model.										
922	LLaMA-2-70B-Int4	35.33GB	wikitext	4-bit quantization model.										
923	LLaMA-2-70B-Int3	26.78GB	wikitext	3-bit quantization model.										

on Baichuan2-13B-Base (Baichuan, 2023) which is an open-source LLM published by Baichuan In-telligence. GLM (Du et al., 2022), developed by Tsinghua Knowledge Engineering Group, is a Gen-eral 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 param-eters 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.

C.2 PROMPTS

<sup>1</sup> Here is a single-answer multiple choice question about Network Implementations. 以下关于网络实现的单选选择题,请直接给出正确答案的选项。
Which TCP/IP routing protocol among the following does not incorporate the subnet mask in its route update messages, thereby hindering its support for subnetting? 以下哪个TCP/IP路由协议在其路由更新消息中不包括子网掩码,从而无法支持子网划分?
A: Routing Information Protocol, version 1 (RIPv1) B: Routing Information Protocol, version 2 (RIPv2) C: Border Gateway Protocol (BGP) D: Open Shortest Path First (OSPF) A: 路由信息协议,版本1 (RIPv1) B: 路由信息协议第二版 (RIPv2) C: 边界网关协议 (BGP) D: 开放最短路径优先 (OSPF)
Let's think step by step. 让我们逐个选项分析:
A: Routing Information Protocol, version 1 (RIPv1) - This is the correct answer B: Routing Information Protocol, version 2 (RIPv2) - RIPv2 does include the subnet mask… C: Border Gateway Protocol (BGP) - BGP is a routing protocol used for large-scale networks,… D: Open Shortest Path First (OSPF) - OSPF supports subnetting and includes the subnet mask… A: 路由信息协议,版本1 (RIPv1) - 正确。RIPv1不包括子网掩码信息,因此支持子网划分。 B: 路由信息协议第二版(RIPv2) - 错误。RIPv2包括子网掩码信息,因此支持子网划分。 C: 边界网关协议(BGP) - 错误。BGP是一种大型互联网路由协议,支持子网划分。 D: 开放最短路径优先(OSPF) - 错误。OSPF是一种内部网关协议(IGP),支持子网划分。
Therefore the answer is : A 因此答案是: A

Figure 13: 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 13. 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 14.

977

978

979

980

981

982

983

984

985

986

987

988

989

990

994

995 996

997

998

999 1000

1001 1002

974 975 976

以下关于网络基础知识的单选选择题,请直接给出正确答案的选项 Which devices can transmit packets across multiple networks and use tables to store network addresses to determine the optimal destination? 什么设备可以在多个网络之间传输数据包,并使用表格存储网络地址以确定最佳目的地? A: Hubs B: Firewalls C: Routers D: Switches A: 集线器 B: 防火墙 C: 路由器 D: 交换机 答: A-集线器……, B-防火墙……, C-路由器……, D-交换机……。所以答案是C。 ... [3-shot examples] ... Here is a single-answer multiple choice question about Network Implementations. 以下关于网络实现的单选选择题,请直接给出正确答案的选项。 Which TCP/IP routing protocol among the following does not incorporate the subnet mask in its route update messages, thereby hindering its support for subnetting? 以下哪个TCP/IP路由协议在其路由更新消息中不包括子网掩码,从而无法支持子网划分? A: Routing Information Protocol, version 1 (RIPv1) B: Routing Information Protocol, version 2 (RIPv2) C: Border Gateway Protocol (BGP) D: Open Shortest Path First (OSPF) A:路由信息协议,版本1(RIPv1) B:路由信息协议第二版(RIPv2) C: 边界网关协议 (BGP) D: 开放最短路径优先 (OSPF) Answer: A-Routing Information Protocol....... So the answer is A. 答:A-路由信息协议·····,所以答案是A。

Here is a single-answer multiple choice question about Networking Fundamentals.

Figure 14: 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.

C.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.

C.4 OVERVIEW PERFORMANCE ON DIFFERENT TEST SETS

Table 8: LLMs' overall performance on wired network operations test set

	English Test Set						Chinese Test Set									
Model	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	/	/	/	/	1	/	/	86.00	1
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 8, Table 9 and Table 10, 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.

1021

1017

#### C.5 PERFORMANCE ON DIFFERENT QUANTIZATION MODELS

1023

Figure 15 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.

vlodel		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
nternLM2-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
LaMA-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
Aistral-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
nternLM2-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
LaMA-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
3aichuan2-13B-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
LaMA-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 la	anguage	for eac	h LLM is i	n <b>bold</b> f	ont.										
r	Fable	10:	LLM	[s' over	rall n	erfoi	man	ce on c	latab	ase o	pera	tions te	est se	t		
					F						r ····			-		
Model				English	Test Se	1						Chinese	Test Se	t		
loder	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
nternLM2-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
LaMA-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
LaMA-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
LaMA-2-7B	22.13	23.74	23.74	26.56	19.32	20.52	28.77	33.60	20.72	20.72	27.16	27.97	21.53	18.51	18.31	17.91
Aistral-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
nternLM2-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
	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
Jemma-2B	of each l	anguage	for eac	h LLM is i	n <b>bold</b> f	ont.										
Gemma-2B Note: The best accuracy of																
Gemma-2B Note: The best accuracy of	or each h	-	ILLE C									Chin	L			
Gemma-2B Note: The best accuracy o		Eng	lish Se	t	1 70			60 -				Chinese S	Set			]
Solution 150-Chat Gemma-2B Note: The best accuracy of		Eng	llish Se 5	t 5.06 54	4.79			60				Chinese S	Set			
Seema-2B Note: The best accuracy of 50 - 55.81 52	2.31	Eng 42.94	llish Se 5	t 5.06 54	4.79	45.55		60 - 50 -				Chinese S	5et 43.05			
Semma-2B Note: The best accuracy of 50 50 50 50 50 50 50 50 50 50 50 50 50	2.31	En <u>c</u> 42.94	ilish Se 5	t 5.06 54	4.79	45.55	1001	60 - 50 - 40 -	39.16	35.49		Chinese S	5et 43.05	37.8	7 35	i.44
1000000000000000000000000000000000000	2.31	Eng 42.94	ilish Se	t 5.06 54	4.79	45.55	racv(%)	60 - 50 - 2 40 -	39.16	35.49		Chinese S	5et 43.05	37.8	7 35	i.44



10

0

zero-shot

1066



few-shot

10

few-shot

zero-shot

LLaMA2-70B-Int4 can achieve an accuracy close to LLaMA-2-70B without quantization. Specif-1072 ically, on English multi-choice questions, the accuracy of the GPTQ model with 4-bit quantiza-1073 tion parameters is 3.50% lower in zero-shot evaluation and 0.27% in few-shot evaluation compared 1074 to LLaMA-2-70B. As for Chinese questions, the accuracy of LLaMA2-70B-Int4 is 3.67% lower 1075 in zero-shot evaluation and 5.18% in few-shot evaluation compared to LLaMA-2-70B. However, 1076 LLaMA2-70B-Int3 has a performance degradation that cannot be ignored. On average, the accuracy 1077 of LLaMA2-70B-Int3 in English set has a 12.46% degradation compared to LLaMA-2-70B and a 1078 9.30% degradation compared to LLaMA2-70B-Int4.



As depicted in Figure 17, the FAE-Score demonstrates a strong positive correlation with Expert Evaluation Score, making it a valuable and effective substitute for automated evaluation.

Original Question	Mock Question
Your network currently utilizes 802.11ac for all client computers. Recently, there has been a relocation of several users from one office space to another, resulting in an increase in the number of users in the area from 20 to approximately 50. As a result, both new and old users have reported experiencing significantly slower network transfer speeds. What is the most probable cause of this issue? A. The current 802.11ac standard is unable to support such a high number of concurrent users. B. The distance between the wireless access point and the users is too great. C. The wireless access noint is unable to accommodate the	Your network uses 802.11ac for all client computers. Recently, several users moved from one office space to another, increasing the users in the area from 20 to about 50. Now, both new and old users are reporting very slow network transfer speeds. What is most likely the cause of the problem? A. 802.11ac can't support that many concurrent users. B. It's too far from the wireless access point.
increased number of users. D. The new users are equipped with 802.11n network cards.	D. The new users all have 802.11n network cards.
L <sub>test</sub> (Model A): 2.126566 L <sub>test</sub> (Model B): 1.665372	L <sub>ref</sub> (Model A): 2.121720 L <sub>ref</sub> (Model B): 2.562153
Δ <i>L</i> (Model A): +0.004846 Δ <i>L</i> (Model B): -0.896781	
Figure 18: An exam	ple for leakage Test.
Figure 18 shows an example for leakage test. N original question, but uses different vocabulary a suggests that there may be potential leakage in the	ote that the mock question is a paraphrase of and phrasing. The difference in $\Delta L$ for Moc e original question.
C.9 CASE STUDY	
2.9.1 FAILURE CASES OF OPSEVAL QUESTIC	5NS
Figure 19 illustrates three instances of failure in ce re typically attributed to either a deficiency in c lomain knowledge.	rtain LLMs. The primary reasons for these fai lomain knowledge or a lack of comprehension
C.9.2 CASE ANALYSIS ON QA METRICS	
igure 20 illustrates an example where BLEU/RC affuenced by the repeated pattern attack, while F.	DUGE fail to pick up important keywords an AE-Score avoids these two issues.
C.9.3 CASE STUDY: WHY ADVANCED SETTIN	NGS SOMETIMES LACK BEHIND
n certain cases, more advanced evaluation metho o understand the potential reasons behind this pl o the guidance provided by the CoT prompts who outputs. Figure 21 is one of the examples where he idea of thinking step by step. Thus, instead of ame to its answer directly. Even though the mo	ods surprisingly lead to poorer results. We and nenomenon. 1) Some models may respond po- en required to think step by step, leading to su CoT failed: the model tested cannot compre- analyzing each option, it repeated the question del correctly answered "FTP server" when a
n English, it failed to give the expected option A prompting when applying the CoT method. 2) Fe hat the task involves generating questions rathe ssues. Figure 22 provides an example to the prob	A. This failed case inspires the need for few ew-shot prompts may lead some models to be or than answering them, resulting in perform elem mentioned above.
C.9.4 CASE STUDY: HOW BAICHUAN OUTPE	RFORMS IN CHINESE
Figure 23 shows an example where Baichuan-13B with correct English analysis from LLaMA-2-12 (3B-Chat itself for comparison. The malfunctions	-Chat failed in the English 3-shot CoT+SC set 3B and correct Chinese analysis from Baich ed output generates an endless analysis for a s

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.





1242 Here is a single-answer multiple choice question about Networking. 1243 以下关于网络的单洗洗择题, 请直接给出正确答案的洗项 1244 What type of server can host files for easy access and downloading, similar to how a web server serves 1945 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? 1246 1247 A: FTP server C: File server D: DNS server B: Proxv server A: FTP服务器 B: 代理服务器 C: 文件服务器 D: DNS服务器 1248 1249 Let's think step by step. 让我们逐个选项分析 1250 1251 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 1253 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 1255 Therefore the answer is : FTP server 1257 因此答案是: C 1259 Figure 21: 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 1261 dark red and the red text are model's corresponding output. 1262 Here are single-choice questions related to networks. Please provide the correct answer 1263 directly 1264 What network segmentation technique reduces the broadcast domain and provides an 1265 added layer of security among users on the same network? 1266 1267 A: VPN B: VLAN C: UPS D: SQL 1268 Zero-shot Answer: 1270 **3-shot Answer:** 1271 A For the following multiple-choice questions about networking, networking, network 1272 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 1274 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 1276 the Internet? A: Virtual private network [Model's output truncated here] 1278 Figure 22: A failed 3-shot example Generated by Qwen-7B-Chat on both zero-shot and few-shot 1279 evaluations under the naive setting in Chinese. 1280 D ANNOTATION GUIDELINE FOR OPSEVAL CATEGORIZATION 1281 1282 D.1 OVERVIEW 1283 1284 In the OpsEval project, we aim to categorize operational and maintenance tasks within the industry. 1285 This categorization process is pivotal for understanding the spectrum of tasks and the required abili-1286 ties to address them effectively. The process involves two primary steps: automated screening using 1287 GPT-4 for initial topic modeling, followed by a manual review process involving domain experts. 1288 1290 D.2 TASK CATEGORIZATION 1291 D.2.1 OBJECTIVE 1293 1294 To categorize questions into one of eight distinct operational tasks based on industry relevance, task 1295 frequency, and significance within operational settings.

297	
98	
99	
0	
1	
2	
3	
4	
5	
6	
7	
8	
9	A user reports that she cannot connect to a server on her network. Ed wants to identify the scope of the
0	problem, so he tries to reproduce the problem on the user's computer. The problem still remains. No
1	other users are reporting this problem. What is the next logical step that Ed should perform to identify the affected area?
2	A. Vorify that the least router is forwarding traffic P. To use forwing the same task as another such as
3	A. Verify that the local router is forwarding traffic. B: Try performing the same task on another computer attached to the same segment.
4	C: Verify that the server is configured properly. D: Verify that the switch the client is connected to is
15	functioning.
6	Let's think step by step. Verify that the Local Router is Forwarding Traffic. This will help determine if
7	across different segments within one VLAN / Subnetwork respectively thereby narrowing down possible
8	areas where issue lies & helping us pinpoint root cause faster without much effort involved here because
9	[Model continues with long sentence and then truncated]
20	Therefore the answer is : [Model doesn't have any output]
21	Let's think step by step.
22	A: Verify that the local router is forwarding traffic - This is not the next logical step. While it's possible that
23	one user.
24	B: Try performing the same task on another computer attached to the same segment - <b>This is a good</b>
	itself rather than the network.
-9	C: Verify that the server is configured properly - This is not the next logical step. If the problem is with the
.0	D: Verify that the switch the client is connected to is functioning - This is a good idea. If the switch is not
. / ) Q	functioning properly, it could cause issues with network connectivity.
0	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
1	服务商也云大败。然而,这个问题限于个人可能定守我问题的唯一原因。 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
20	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。
2	U: 短证服务
	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
5	到各尸峏与版 <b>穷</b> 舔的通信。虫然这个迹坝有起米很有追埋,但在很多情况下,这个问题开个一定是田士 交换机故障导致的。
0	因此答案是: B
	Figure 23: A failed English answering example Generated by Baichuan 12D Chat on faw shot
	agure 23. A raneu English-answering example Generated by Datchuan-13D-Unal on Iew-shot
	auon under me COTTOC setting in ootn English and Chillese.
9	
U G	
1	

1350	D.2.2	STEPS
1352	1.	Review Initial Categorization: Begin with the insights provided by GPT-4's topic model-
1353	2.	<b>Understand Task Definitions</b> : Familiarize yourself with the details of the eight distinct
1354		tasks outlined in the provided Appendix. Each task has specific criteria and examples to
1356	3.	guide your categorization. Assign Tasks: For each question, decide which of the eight tasks it belongs to. A question
1357	01	should be categorized based on its core focus and the operational activity it pertains to.
1358 1359	4.	<b>Justification</b> : 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
1360		than one category. Use the task demittions as a guide to support your decision.
1361 1362	D.2.3	DETAILED TASK CATEGORIZATIONS
1363	1.	General Knowledge: Questions related to foundational concepts and practices in the Ops
1364 1365	2	domain. Fault Analysis and Diagnostics: Questions focusing on identifying and solving discrep-
1366	2.	ancies or faults in systems or networks.
1367	3.	<b>Network Configuration</b> : Questions about optimal configurations for network devices to
1368	4.	<b>Software Deployment</b> : Questions dealing with the distribution and management of soft-
1370	5	ware applications.
1371	5.	and setting up alert mechanisms.
1372 1373	6.	Performance Optimization: Questions aimed at enhancing network and system perfor-
1374	7.	Automation Scripts: Questions involving the creation of scripts to automate processes and
1375	0	reduce manual intervention.
1376 1377	8.	multiple categories.
1378		
1379 1380	D.2.4	TASK CATEGORIZATION TEMPLATE
1381	Quest	ion ID:
1383	Assig	ned Task:
1385	Justif	cation: [Provide a brief explanation for the task assignment here]
1386		
1387	D.2.5	EXAMPLE FOR TASK CATEGORIZATION
1389		
1390	Quest	ion ID: 001
1391	Quest	allowing legitimate traffic?
1393	Assig	ned Task: Network Configuration
1394	Justifi work	device (firewall) to ensure security and efficient operation, aligning perfectly with the
1395	'Netw	ork Configuration' task.
1397		
1398 1399		
1400	ы.5 Р	ADILITI CALEGORIZATION
1401	D.3.1	OBJECTIVE
1402	To class Analytic	sify questions based on the required cognitive ability to answer them: Knowledge Recall, cal Thinking, or Practical Application.

D.3.2	STEPS					
1.	<b>Review Definitions</b> : Read the descriptions of the three abilities in the provided Appendix. Each ability category has distinct characteristics and examples.					
2.	<b>Evaluate Questions:</b> 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.	<b>Assign Ability Level</b> : 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.	<b>Justification</b> : 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 catego- rization.					
D.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.	<b>Practical Application</b> : Involves applying knowledge or analytical insights to provide actionable recommendations.					
D.3.4	Ability Categorization Template					
<ul> <li>Question ID:</li> <li>Question: [Insert question text here]</li> <li>Assigned Ability:</li> <li>Justification: [Provide a brief explanation for the ability level assignment here]</li> </ul>						
D.3.5	EXAMPLE FOR ABILITY CATEGORIZATION					
Quest ing fr Assig edge tions,	ion ID: 002 Question: How would you optimize the performance of a network experienc- equent bottlenecks? ned Ability: Practical Application Justification: The question requires applying knowl- of network systems and performance optimization techniques to propose specific solu- hence it falls under 'Practical Application'.					
D.4 C	eneral Guidelines					
•	<b>Consistency</b> : 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.					
•	<b>Collaboration</b> : When in doubt, discuss challenging questions with fellow experts. Collaboration can help clarify ambiguities and refine the categorization process.					
•	<b>Documentation</b> : Keep detailed notes on your decisions, especially for questions that re- quired significant deliberation. This documentation will be valuable for future reference and analysis.					
By follo of opera the oper	wing these guidelines, you will contribute to a comprehensive and nuanced categorization tional tasks and required abilities. This effort is crucial for enhancing our understanding of ational landscape and the diverse skills professionals need to navigate it effectively.					
	D.3.2 1. 2. 3. 4. D.3.3 1. 2. 3. D.3.4 Quest Quest Assig Justifi D.3.5 Quest ing from Assig edge of tions, D.4 G • • • • • • • • • • • • •					