
CLOUDEVAL-YAML: A PRACTICAL BENCHMARK FOR CLOUD CONFIGURATION GENERATION

Yifei Xu^{*1,3}, Yuning Chen^{*1,4}, Xumiao Zhang^{*2}, Xianshang Lin¹, Pan Hu^{†1},
Yunfei Ma¹, Songwu Lu³, Wan Du⁴, Z. Morley Mao², Dennis Cai¹, Ennan Zhai¹
¹Alibaba Cloud ²University of Michigan ³UCLA ⁴UC Merced

*Co-first authors. †Corresponding author.

ABSTRACT

Among the thriving ecosystem of cloud computing and the proliferation of Large Language Model (LLM)-based code generation tools, there is a lack of benchmarking for code generation in cloud-native applications. In response to this need, we present `CloudEval-YAML`, a practical benchmark for cloud configuration generation. `CloudEval-YAML` tackles the diversity challenge by focusing on YAML, the de facto standard of numerous cloud-native tools. We develop the `CloudEval-YAML` benchmark with practicality in mind: the dataset consists of hand-written problems with unit tests targeting practical scenarios. We further enhance the dataset to meet practical needs by rephrasing questions in a concise, abbreviated, and bilingual manner. The dataset consists of 1011 problems that take more than 1200 human hours to complete. To improve practicality during evaluation, we build a scalable evaluation platform for `CloudEval-YAML` that achieves a 20 times speedup over a single machine. To the best of our knowledge, the `CloudEval-YAML` dataset is the first hand-written dataset targeting cloud-native applications. We present an in-depth evaluation of 13 LLMs, leading to a deeper understanding of the problems and LLMs, as well as effective methods to improve task performance and reduce cost. We release the dataset along with the evaluation framework at <https://github.com/alibaba/CloudEval-YAML>.

1 INTRODUCTION

Cloud computing is a \$545 billion business (clo, 2023) with a thriving ecosystem of open-source, cloud-native tools and applications. There are more than 1000 cloud-native applications in the Cloud Native Computing Foundation landscape (cnc, 2023) that have attracted over 4 million Github stars. Many businesses, from small startups to large corporations, rely on cloud services for their operations. On the other hand, recent LLM-based code generation tools such as ChatGPT (cha, 2023; gpt, 2023), Codex (cod, 2023), Copilot (cop, 2023), and Llama (Touvron et al., 2023a;b) have been shown to significantly increase the efficiency of software developers (Vaithilingam et al., 2022; Kabir et al., 2023).

Despite the significance of cloud computing and the effectiveness of LLM-based code generation tools, there is a lack of benchmarking for code generation in cloud-native applications. Understanding the performance of various LLM models in the context of cloud applications remains an open question. It is crucial to establish benchmarks that facilitate direct comparisons and foster the evolution of these models and methods.

We attribute the absence of such benchmarks to the following challenges: First, cloud applications are diverse in terms of programming languages and application interfaces,

which complicates the construction of a unified benchmark. Second, cloud applications are usually deployed on a large scale, which makes it hard to evaluate the generated code.

In this paper, we present `CloudEval-YAML`, a first-of-its-kind benchmark for cloud configuration generation. It tackles the aforementioned diversity challenge by focusing on YAML, the de facto standard of numerous cloud-native applications. It is both human and machine-readable, thus simplifying the configuration of cloud applications. 90 out of the top 100 most-starred cloud-native applications on GitHub use more than 10 YAML files. We include more detailed statistics in Appendix A.

The `CloudEval-YAML` is designed around practicality, featuring a dataset that consists of 1011 problems and a comprehensive, scalable benchmark framework. The dataset targets realistic problems from a wide range of sources that take over 1200 human hours to complete. Each problem in the dataset includes a question context composed of natural language descriptions and an optional code segment. Additionally, there is a labeled reference YAML file and a unit test script in each problem for evaluation purposes. Figure 1 shows an illustrative example of a problem.

We further enhance the dataset with practical data augmentation to meet the needs of actual users consisting of simplification and translation rewriting. The simplification in-

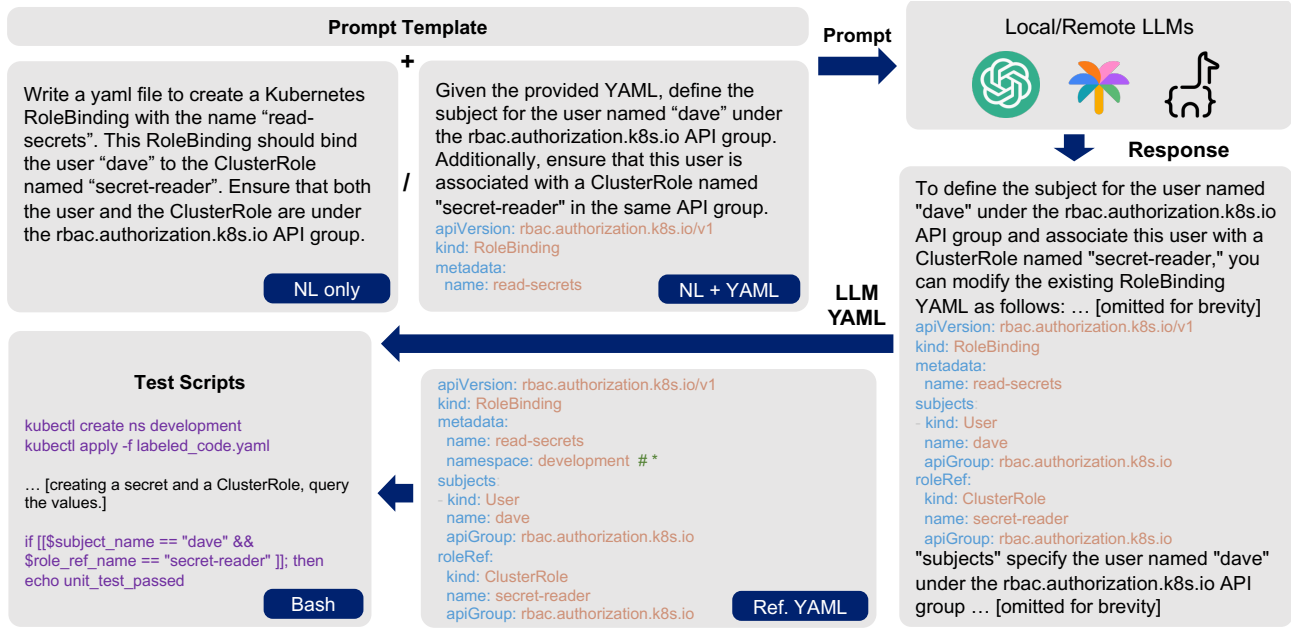


Figure 1: An example problem of the CloudEval-YAML dataset, including a problem specification in natural language with an optional sample YAML file as prompt to LLMs, a reference YAML file, and a bash unit test script, to evaluate the YAML output from the LLM.

involves using concise language with many abbreviations, while translation uses the native language of Chinese used by our cloud operation teams. We integrate GPT-4 into the rewriting pipeline to speed up the process, yet we manually review each of them to ensure data quality.

The CloudEval-YAML benchmark also includes a scalable automated evaluation platform with unit tests to ensure functional correctness, as well as a distributed evaluation cluster to score the generated code efficiently. We employ a cloud-based evaluation framework coupled with shared docker image caching. With a cluster of 64 4-core 8GB machines, we managed to complete the evaluation of 1011 problems in 30 minutes, which would take 10 hours on an individual machine.

We perform an in-depth evaluation on several popular open-source/proprietary language and code generation models with CloudEval-YAML, including Llama/Llama2 (Touvron et al., 2023a;b), Wizardcoder (Luo et al., 2023), Google PaLM-2 (Anil et al., 2023; Chowdhery et al., 2022), and OpenAI GPT-3.5/GPT-4 (Cha, 2023; gpt, 2023). The evaluation includes a comprehensive benchmark of different evaluation metrics, a performance analysis from various perspectives, and experiments with methods to improve generation performance such as multi-sample generation and few-shot prompting. Lastly, we train models that predict unit test results to reduce cost.

The benchmark leads to a series of notable observations, including: 1) Proprietary models outperform open-source models by a large gap. 2) Dedicated code generation models

perform poorly on this task. 3) Length of answer and application category are major factors that affect performance. 4) Even the best-performing model makes simple mistakes. 5) Multi-sample generation can provide remarkable gains and can make cheaper models more cost-effective. 6) Few-shot prompts do not show a significant gain on this task. 7) Predicting unit test scores using other metrics can provide an estimation of the ranking.

In summary, our contributions are as follows:

- We present CloudEval-YAML benchmark that includes the first hand-written dataset with 1011 practical problems for cloud-scale applications.
- We present the design of a scalable, automated evaluation platform consisting of a computing cluster to evaluate the generated code efficiently for various performance metrics.
- We present an in-depth evaluation of 13 LLMs with CloudEval-YAML, leading to a deeper understanding of the performance of LLMs in the context of cloud configuration, as well as effective methods to improve task performance and reduce cost.

The paper is organized as follows: We first introduce the CloudEval-YAML dataset in § 2. Then, we detail the CloudEval-YAML benchmark platform as well as the evaluation criteria in § 3, followed by § 4 where we present the experiments and results. We discuss related works in § 5 and conclude the paper in § 6.

2 THE CLOUDEVAL-YAML DATASET

2.1 Overall Structure

The overall structure of CloudEval-YAML dataset, illustrated with a simple example, is shown in Figure 1. It consists of the following components:

- **Prompt template:** The prompt template is added to the beginning of each problem to provide context for the model, as well as to specify the output format of the desired answer. We use the same template for all problems, which can be found in Appendix B.
- **Natural Language (NL) problem description and optional YAML context:** Within the *CloudEval-YAML* dataset, problem descriptions come in two forms: those consisting of natural language only and those paired with a YAML context. The YAML context serves as an input to guide the LLMs in infilling, modifying, or enhancing the functionality of the YAML context.
- **Reference YAML with label:** The dataset also includes a reference YAML file which serves two purposes: First, it can be used as a reference to evaluate the generated YAML file. For example, one can calculate text-level similarities using metrics such as BLEU. However, we find that text-level metrics often fail to consider important properties of YAML such as the fact that object order is not critical. Such oversights can significantly affect the BLEU score. We describe how we tailor the metrics for YAML in §3. To facilitate YAML-aware comparisons, we include three kinds of labels as comments in the YAML file: 1) wildcard match (labeled with #*); 2) exact match (by default, no labeling required); and 3) conditional match (e.g., #v in [2, 3, 4]). Aside from text-level and YAML-aware metrics, we also use the reference YAML file to facilitate the development and verification of the unit test script.
- **Test script:** To benchmark the functional correctness of the generated YAML file, we write automated test scripts to set up the environment and validate the functionality using assertions. For example, if all checkpoints pass, the script outputs `unit_test_passed` to a log file to be further processed and aggregated by the scoring script. The test script also includes a clean-up function ensuring the environment is reset after each test, so the next test can start with a clean environment.

To ensure the authenticity and practicality of the dataset, we collect problems from carefully selected sources including official documentation websites, popular issues from StackOverflow, and highly-ranked blog posts. The problems are hand-picked from these sources based on the following guidelines: 1) The problem should be clearly defined and its purpose should be easily understandable. We avoid problems that are ambiguous in intent or heavily dependent on other YAML files. 2) We collect a diverse set of problems

from common cloud applications to ensure our benchmark is comprehensive. We also pick problems of varying difficulty levels, as well as different types of problems that cover common uses of cloud-native tools, as indicated in Table 2. 3) While it is possible to include scripts like bash, python, or perl in YAML files, we focus on testing the LLM’s capability to generate YAML files for cloud applications, rather than creating specific scripts, and avoid problems that require the generation of complex scripts. After hand-picking a problem, we write the problem description and unit tests, then label the YAML file. Part of unit tests are hand-adapted from online sources. More samples of the dataset are shown in Appendix C for an in-depth view of the dataset.

2.2 Practical Data Augmentation

We conduct a survey with the operation team of a leading cloud provider to align the dataset to representative users. The survey leads to two findings: firstly, actual users often use domain-specific abbreviations and concise language instead of full sentences to express their intentions; secondly, a significant number of users opt to use their native language rather than English. Based on these findings, we apply data augmentation on the original questions to make the dataset more practical, including simplification and translation. We illustrate the process with an example as shown in Figure 2.

We use GPT-4, a state-of-the-art LLM, to assist the data augmentation process: ① we start by writing three examples for simplification and translation. For simplification, we use short and clear language with as many abbreviations as possible. For translation, we use the daily language of Chinese used by our cloud operation teams. The full prompts for both tasks can be found in Appendix D. ② However, we cannot fully depend on GPT-4 as it can make mistakes. So we manually review the simplified and translated question drafts created by GPT-4 to make sure they keep their original meaning and complexity, ensuring experts can answer them. For the translated drafts, we employ native speakers with cloud development experience to ensure the translated version matches the standard practices of the target developers.

Statistics of the original and augmented datasets are displayed in Table 1. We reduce the average word/token count by 25.7%/20.9% respectively via simplification. The translated dataset uses fewer words/tokens, but it is not comparable to the original dataset due to language discrepancy.

Table 1: Statistics of Practical Data Augmentation

	Original	Simplified	Translated
Count	337	337	337
Avg. words	99.40	73.86 (-25.7%)	57.18
Avg. tokens	508.9	402.5 (-20.9%)	378.5

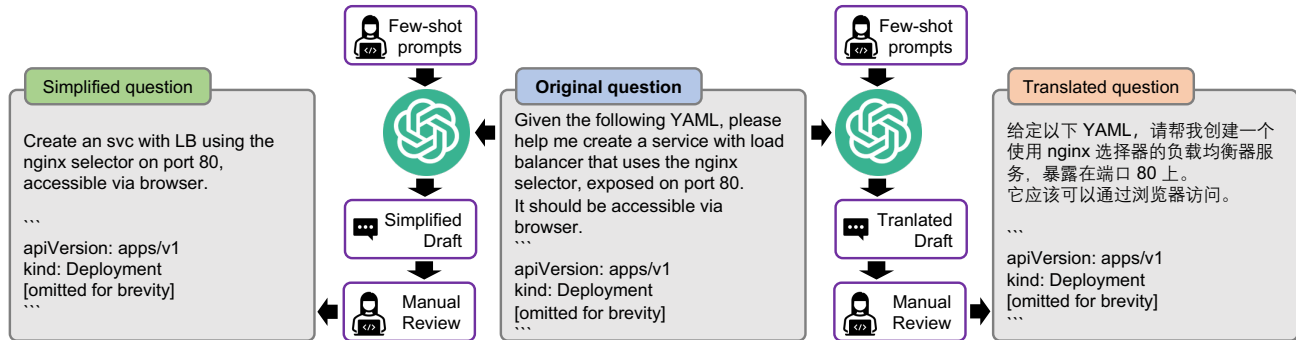


Figure 2: The practical data augmentation framework with examples of a simplified and a translated question.

Table 2: Statistics of CloudEval-YAML dataset

Statistics	Kubernetes						Envoy	Istio	Total / Avg. / Max
	pod	daemonset	service	job	deployment	others			
Total Problem Count	48	55	20	19	19	122	41	13	337
Avg. Question Words	77.06	80.91	71.35	73.74	94.84	69.48	275.56	73.00	99.40
Avg. Lines of Solution	18.67	23.58	15.00	20.37	29.00	19.74	85.85	14.92	28.35
Avg. Tokens of Solution	64.02	71.91	41.40	74.53	79.42	58.78	242.34	39.54	84.28
Max Tokens of Solution	150	111	83	163	140	194	531	53	531
Avg. Lines of Unit Test	8.52	8.58	11.25	7.68	12.53	17.74	11.56	20.00	13.14

2.3 Statistics of the CloudEval-YAML dataset

We invest more than 1200 human hours to create CloudEval-YAML dataset. The statistics of the dataset are shown in Table 2. The dataset covers a variety of functionalities in Kubernetes, such as `pod`, `daemonset`, `service`, `job`, and `deployment`, while also providing insights into other software like Envoy and Istio, offering a comprehensive overview of each system software’s major features and applications. Aside from its wide coverage of real-world applications, the dataset also includes practical problems that are more challenging than other hand-written datasets such as HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) in terms of problem/solution length. For example, the average line of solution in this dataset is 28.35, which is $4\times$ as HumanEval (6.3) and MBPP (6.7). We take extra caution in creating the dataset, including limiting the maximum token count (531) of solutions so it would not exceed the context length of small models, as well as detailed unit tests (13.14 lines on average) to ensure the functional strictness of the test scripts.

3 THE CLOUDEVAL-YAML BENCHMARK PLATFORM

The overall framework of CloudEval-YAML is depicted in Figure 3. It is written in 6.3k lines of Python and Bash code, consisting of three major components: YAML gen-

eration, cloud evaluation, and score calculation. We first introduce the design of YAML generation and score calculation, and then explain the design of cloud evaluation given the complexity of score calculation.

3.1 YAML Answer Generation

As explained in §2, each problem in our dataset includes a file with a problem description. We create prompts for LLMs by combining the prompt template with the problem. These prompts are then processed by the query module to generate a YAML file.

The query module serves as a universal interface for various local and remote models. There are several reasons for doing so. First, it simplifies the differences between different local and remote APIs, offering a consistent interface. This makes it easier to add new models, which only requires adding a query wrapper and a model initializer. Second, it optimizes throughput by using parallel processing to fully utilize the rate-limit and auto-scaling features of Model-as-a-service providers¹. As an example, it is straightforward to start 128 raylets with the help of `ray` (Moritz et al., 2018) in the query module, which can significantly increase the speed by two orders of magnitude. For local models, the module

¹For example, `replicate.com` provides autoscaling by running multiple copies of a model for no additional cost, as users only pay for the time that *actually running the inference*, rather than idling/booting time.

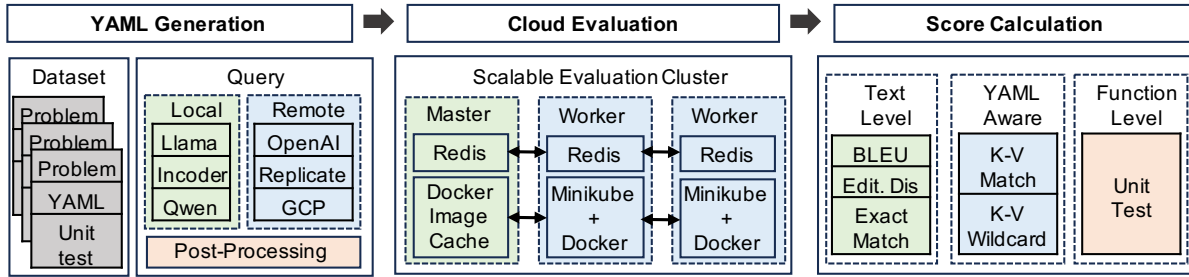


Figure 3: Workflow of the CloudEval-YAML benchmark platform.

automatically checks the available GPU memory and adjusts the batch size to speed up inference.

Post-processing: Although we explicitly require LLMs to answer with YAML only, the response often contains text descriptions wrapped around a valid YAML file. We apply the following post-processing policies to extract clean YAML files from such responses:

- Remove all content before the line with the keyword `Here`, as it is commonly found before the YAML file in responses from several LLMs.
- Remove all content before the line with the keyword `apiVersion:` (for Kubernetes) or `static_resources:` (for Envoy) since they typically mark the start of a YAML file.
- Extract text enclosed by the following delimiters: `````, `<code>` and `</code>`, `\begin{code}` and `\end{code}`, `START SOLUTION` and `END SOLUTION`.

3.2 Performance Score Calculation

We calculate comprehensive scores using three distinct methods that include a total of 6 metrics to cover different aspects of the generated YAML files. The first method, known as the text-level score, uses metrics such as BLEU, Edit Distance, and Exact Match. The second method, referred to as the YAML-aware score, uses the Key-Value Exact Match and Key-Value Wildcard Match. The third method, the function-level score, uses Unit Tests. Here is how we calculate each:

- **BLEU:** Bilingual Evaluation Understudy (Papineni et al., 2002) is a common metric used to evaluate the quality of machine-generated translations. It measures the similarity between the generated YAML and the reference YAML. We use the standard implementation from the NLTK (Bird et al., 2009). The BLEU score ranges from 0 to 1, with the higher scores being more desirable.
- **Edit Distance:** In some scenarios even imperfect configuration could still be useful if users can fix the error by modifying a few words. The Edit Distance metric is calculated by comparing the number of lines to edit between the generated YAML and the reference YAML

using Python standard library `difflib.Differ`. We scale the edit distance by the size of the reference YAML using $1 - \text{edit_distance} / \text{len}(\text{reference_YAML})$. As a result, the edit distance score ranges from 0 to 1; the higher, the better.

- **Exact Match:** Opposite to edit distance, the exact match score is a very strict metric that requires the generated YAML to be exactly the same as the reference YAML. The output is either 0 (not match) or 1 (exact match).
- **Key-Value Exact Match:** Different from the exact match that ignores the fact that order does not matter in YAML, key-value exact match loads both the generated and reference YAML files into dictionaries and checks if the resulting dictionaries are the same or not, so the output is either 0 (not match) or 1 (exact match).
- **Key-Value Wildcard Match:** Similar to the key-value exact match, we also load both YAML files into the dictionary. However, with the help of labeling in the reference YAML file, we can tell what matters and what is not critical. For example, sometimes it is acceptable to start a cluster with `image: ubuntu:22.04` or `image: ubuntu:20.04`, so the label in reference YAML could be `image: ubuntu:22.04 # v in ['20.04', '22.04']` and either version will be considered correct. We implement this key-value wildcard match using a tree with leaf nodes marked in `exact/set/wildcard` match and then calculate the IoU (intersection over union) of dictionaries from the generated and reference YAML. The score ranges from 0 to 1; the higher, the better.
- **Unit Test:** CloudEval-YAML ensures the functional correctness of the generated YAML files by running unit test scripts crafted by domain experts. For Kubernetes-focused applications, including Kubernetes and Istio, Minikube (min, 2023) offers the capability to set up virtual Kubernetes clusters within a local testing environment. The `kubect1` command set, the standard tool for managing actual Kubernetes clusters, functions identically on these virtual clusters. It is used for tasks including setting up the environment, applying the YAML files, and monitoring the status. We use Docker to establish the cluster and perform testing on containers directly

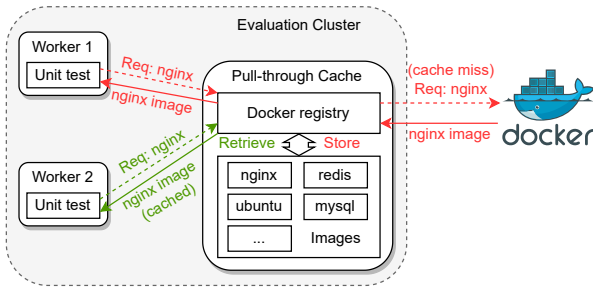


Figure 4: Architecture of shared Docker image caching.

for Envoy applications. A unit test is formulated for each problem, tailored to the expected functionality, and it outputs 1 if successful and 0 if not.

3.3 Cloud-based Evaluation Framework

In contrast to the text-level and YAML-aware scores that take only 21.9 seconds to compute over the entire dataset, running the unit tests in the real environment is very time-consuming: it usually takes several minutes to create the cluster, pull corresponding images, initialize and apply configurations, as well as clear up the environment. This sequential process, especially when scaled to hundreds of problems, becomes a significant bottleneck as it can take more than 10 hours to complete on a single machine. To expedite the evaluation, we employ the following techniques:

Scalable Evaluation Cluster: We design a scalable evaluation cluster to serve as the unit testing backend. Distinct from the previously mentioned Kubernetes and Docker clusters that run on a local machine, this cluster consists of a coordinator node and multiple worker nodes that span several virtual machines. Central to this system, the coordinator employs a Redis database to manage unit test contexts, inputs, and outputs associated with each problem and benchmark user. Users can dispatch their unit testing jobs to the coordinator, which are claimed by available worker nodes. Subsequently, these workers relay the results back to the coordinator. This paradigm enables automatic parallelization of unit testing, while also ensuring that users can easily monitor progress and access results. Additionally, the distributed design of the evaluation cluster allows for dynamic scaling as needed.

Shared Docker Image Caching: Even with the aforementioned scalable cluster, pulling Docker images from the Dockerhub repeatedly not only consumes excessive Internet bandwidth but may also result in an unexpected timeout due to bandwidth variations. To solve this problem, we set up a local docker hub registry that serves as a pull-through cache on the coordinator node, as shown in Figure 4. This allows workers to share the cached images, removing the need to pull the same image from the Internet if another worker has already downloaded it.

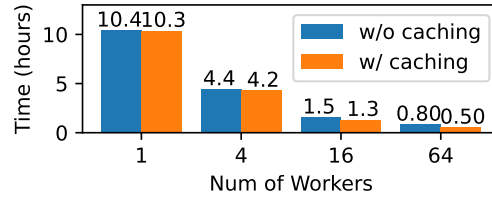


Figure 5: Evaluation time over all 1011 problems.

Micro-benchmark of the cloud-based evaluation time:

We run a micro-benchmark on the evaluation time with our cloud-based evaluation framework. The results are shown in Figure 5. We provide the system with 100 Mbps overall bandwidth for Internet access. With a cluster consisting of 64 workers, each equipped with 4 CPU cores and 8GB memory, and shared docker image caching enabled, the evaluation of all 1011 problems is completed in less than 30 minutes. Compared to a single machine, which needs over 10 hours for the same task, this is over 20× speed improvement, with 13× from parallel unit testing and 1.6× from shared Docker image caching. This enables the scalability across multi-GPU (Feng & Jeon, 2023) and big memory systems (Feng et al., 2023).

3.4 Running Cost of the Benchmark

The main expenses for evaluation are running LLM inference and unit tests in the cloud. Table 3 provides a sample cost breakdown of evaluating all 1011 problems for two different inference options (GPT-3.5 API and Llama-7b over replicate.com) and 3 cloud evaluation settings (launching a GCP cluster with 64 standard 4-core 8GB instances, using spot instances to save costs, and using only a single spot instance). The most cost-effective method is using the GPT-3.5 API and a single GCP spot instance, which costs only \$1.31 per run.

Table 3: Sample Running Cost of the Benchmark in \$.

LLM Inference	Cloud Evaluation	Total Cost
- GPT-3.5: \$0.60	- GCP spot ×1: \$0.71	\$1.31
- Llama-7b: \$2.90	- GCP spot ×64: \$2.20	~
	- GCP std. ×64: \$5.51	\$8.41

4 EVALUATIONS ON CLOUDEVAL-YAML

We present an in-depth evaluation of 13 LLMs on the CloudEval-YAML in this section, including a comprehensive benchmark of all evaluation metrics and performance analysis from various perspectives including problem/dataset types and failure modes; we also experiment with methods to improve generation performance such as multi-sample generation and few-shot prompting. Lastly, we train models that predict unit test results to reduce cost.

4.1 Comprehensive Benchmarks and Analysis

We conduct benchmarks on both open-source and proprietary models of varying sizes, including GPT-3.5 (cha, 2023), GPT-4 (gpt, 2023), GPT-4 Turbo (gpt, 2023), PaLM 2 (pal, 2023), Llama (Touvron et al., 2023a), Llama 2 (Touvron et al., 2023b), Code Llama (Rozière et al., 2023), and WizardCoder (Luo et al., 2023). We focus on the "chat" or "instruct" version of models to better fit in the Q/A nature of the configuration generation. The results of our comprehensive benchmark are presented in Table 4. GPT-4 Turbo performs the best in all text-level, YAML-aware and function-level scores. The unit test score is 0.567, meaning that it could pass more than 56% of the unit tests. Among all the open-source models, Llama-2-70b-chat is the winner that performs best in all metrics other than the `exact_match`.

According to the results in Table 4, we conclude: ❶ Proprietary models such as GPT-3.5 and GPT-4 are way ahead across all metrics, and the gap between them and the best performing open-source models is larger than in similar benchmarks like HumanEval (Chen et al., 2021). On HumanEval, Llama 2 (70B) is able to achieve a score of 29.9 compared to 48.1 and 67.0 ($1.61\times$ and $2.24\times$ as Llama 2 (70B)) for GPT-3.5 and GPT-4 respectively. On the unit test score of CloudEval-YAML, llama-2-70b-chat scores 0.085 whereas GPT-3.5 and GPT-4 score 0.412 and 0.515, which is $4.84\times$ and $6.06\times$ as that of llama-2-70b-chat respectively. ❷ Code LLMs typically perform poorly on CloudEval-YAML compared to general LLMs with similar or even smaller sizes in terms of the Unit Test score: wizardcoder-34b-v1.0 scores 0.056, while llama-2-13b-chat gets a higher score of 0.067 with less than half the model size. It may be related to the dataset used in the fine-tuning process. However, further investigation is required to uncover the underlying reasons.

Performance analysis across various problem types: The analysis of test results on different types of questions is shown in Figure 6 (see Table 9 in Appendix E for more details). From the analysis, we draw the following conclusions: ❶ While it is not surprising that problems with longer answers are more challenging, we observed a steep decrease in scores for medium-long ([15, 30]) problems. Whether it is due to the emergence of larger models worth future investigation. ❷ The negative correlation between scores and the length of questions is relatively weaker than the length of answers. The longer inputs may indicate more complex requirements, but may also help to provide more context to narrow down the search. ❸ The presence of a code context does not have a substantial influence on the performance most of the time. However, we notice that models with an index from 7 to 10 perform better with code context, indicating that the code may provide some clue about the

answer. ❹ Envoy questions are clearly more challenging than Kubernetes and Envoy, which is not surprising due to their longer code lengths compared to other applications as shown in Table 2.

Performance analysis of the practically augmented datasets: As mentioned in §2.2, the CloudEval-YAML dataset comprises both original and simplified/translated questions, tailored to fulfill practical requirements. We compare the performance of these questions, and the result is shown in Table 5. We focus on unit test scores only for simplicity. The results suggest that: ❶ Simplification of problems generally leads to lower performance, but it does not affect large models (such as GPT-4 changed from 179 to 164) as much as small models (Llama-2-13b-chat: 26 to 17); ❷ Code-specific (wizardcoder-34b-v1.0: 24 to 2) and small (Llama-2-7b-chat: 13 to 5) models are severely affected by translation, while large models keep up their performance relatively well (Llama-2-70b-chat: 30 to 32), suggesting the need to fine-tuning small/code-specific models on the bilingual dataset.

Failure Mode Analysis: Not all failures are the same: some are close to the correct answer, while others are completely wrong. By analyzing common failure modes, we can understand the weaknesses of each model and find methods to improve the performance. We group the answers into 6 categories, sorted by how close they are to the correct answer: 1) empty or less than 3 lines; 2) longer than 3 lines but does not contain the `kind` field²; 3) contains `kind` but not a complete YAML file; 4) valid YAML but `kind` field is incorrect; 5) valid YAML, `kind` field is correct but the unit test fails; 6) correct YAML that passes the unit test. We analyze the failure modes of the original dataset and the statistics of the result are shown in Figure 7.

An interesting fact we observed from Figure 7 is that ❶ the best-performing model GPT-4 makes more Category 1 errors (i.e. simple mistakes) than both Llama2-7B/70B. But given that such errors could be easily filtered, we expect that the performance of GPT-4 could be further improved by implementing a basic format check to filter out such errors and regenerate new ones. On the other hand, ❷ both Llama2-7B/70B make a large number of Category 5 errors, suggesting that they are able to get the general idea most of the time, but are not accurate enough to pass unit tests.

4.2 Multi-sample Generation

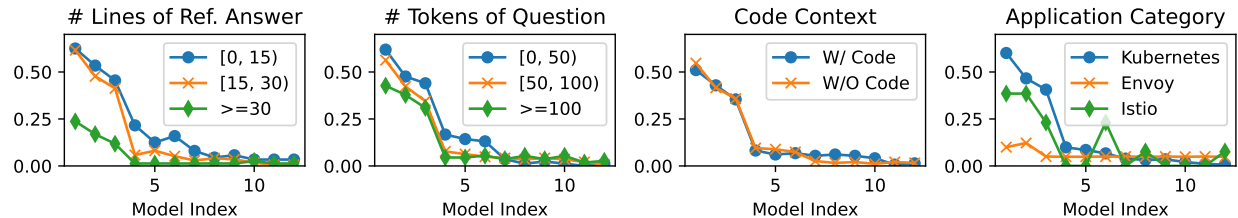
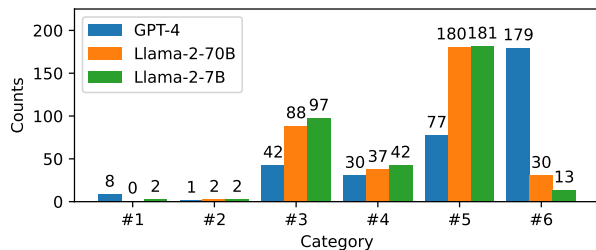
Sometimes users may not be satisfied with the generated result and might generate additional samples to select the most suitable one. To evaluate the performance of mod-

²The field `kind` exists in most Kubernetes configurations. We search for `static_resources` field instead for Envoy configurations.

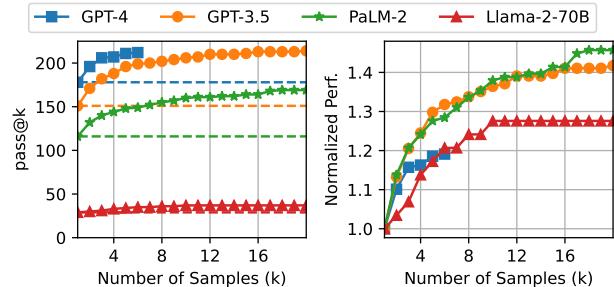
Table 4: Average score of zero-shot benchmark on different models (the higher the better)

Ranking	Model			Text-level Score			YAML-Aware Score		Function-level Score
	Name	Size	Open Source	BLEU	Edit Dist.	Exact Match	Key-value Exact	Key-value Wildcard	Unit Test ↓
1	GPT-4 Turbo	?	N	0.649	0.551	0.099	0.208	0.667	0.561
2	GPT-4	?	N	0.629	0.538	0.092	0.198	0.641	0.515
3	GPT-3.5	?	N	0.612	0.511	0.075	0.154	0.601	0.412
4	PaLM-2-bison ¹	?	N	0.537	0.432	0.040	0.092	0.506	0.322
5	Llama-2-70b-chat	70B	Y	0.355	0.305	0.000	0.020	0.276	0.085
6	Llama-2-13b-chat	13B	Y	0.341	0.298	0.000	0.016	0.265	0.067
7	Wizardcoder-34b-v1.0	34B	Y	0.238	0.247	0.007	0.013	0.230	0.056
8	Llama-2-7b-chat	7B	Y	0.289	0.231	0.000	0.009	0.177	0.027
9	Wizardcoder-15b-v1.0	15B	Y	0.217	0.255	0.002	0.002	0.226	0.026
10	Llama-7b	7B	Y	0.106	0.058	0.004	0.005	0.069	0.023
11	Llama-13b-lora	13B	Y	0.101	0.054	0.001	0.003	0.065	0.021
12	Codellama-7b-instruct	7B	Y	0.154	0.174	0.001	0.001	0.124	0.015
13	Codellama-13b-instruct	13B	Y	0.179	0.206	0.002	0.002	0.142	0.012

¹ The PaLM API supports English only at the time of submission so we averaged the score excluding translated questions.


Figure 6: Performance analysis from four perspectives. The x-axis is the model index according to rank in Table 4.

Figure 7: Failure analysis of different models, grouped in 6 modes.

els in such a scenario, we generate multiple samples and evaluate the performance over the original dataset. We select the best-performing open/closed-source models including Llama-2-70B, PaLM-2, GPT-3.5 and GPT-4 to evaluate their performance with multi-sample generation. We leave parameters that control the randomness of the output to default for proprietary models and set the values of Llama-2-70B to 0.75/0.9/50 for temperature, top_p and top_k respectively. We define pass@k as the number of passed problems where a problem is considered passed if any of its k samples passes the unit test (Kulal et al., 2019).


Figure 8: Pass@k scores of 4 different models in CloudEval-YAML.

The result is shown in Figure 8³. The normalized performance indicates an improvement compared with 1-sample results. We observe that: ① 20-sample generation could improve the unit test score of Llama-2-70B/PaLM-2/GPT-3.5 by 28%/47%/42% respectively, indicating that multi-sample generation could be a good choice to improve the performance if there is a unit test or the user can manually select the best result. ② Although the curves of different models will not cross over each other, it is possible to achieve

³We run GPT-4 for only 6 samples due to the API rate limit.

Table 5: Number of problems passing unit test on the original and practically augmented datasets

Model	Data Set		
	Original	Simplified	Translated
Name			
GPT-4	179	164 (-15)	178 (-1)
GPT-3.5	142	143 (+1)	132 (-10)
PaLM-2-bison	120	97 (-23)	N/A ¹
Llama-2-70b-chat	30	24 (-6)	32 (+2)
Llama-2-13b-chat	26	17 (-9)	25 (-1)
Wizardcoder-34b-v1.0	24	31 (+7)	2 (-22)
Llama-2-7b-chat	13	9 (-4)	5 (-8)
Wizardcoder-15b-v1.0	12	11 (-1)	3 (-9)
Llama-7b	12	7 (-5)	4 (-8)
Llama-13b-lora	8	9 (+1)	4 (-4)
Codellama-7b-instruct	5	6 (+1)	4 (-1)
Codellama-13b-instruct	5	2 (-3)	5 (+0)

¹ The PaLM API supports English only.

Table 6: Unit test score on few-shot prompting

Model	Number of Shots			
	0-shot	1-shot	2-shot	3-shot
Name				
GPT-3.5	142	150 (+8)	143 (+1)	154 (+12)
Llama-2-70b-chat	30	23 (-7)	26 (-4)	29 (-1)
Llama-2-7b-chat	13	14 (+1)	13 (+0)	15 (+2)

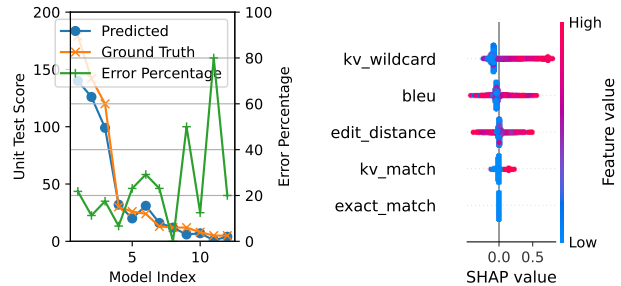
the same or even better performance with multiple samples. For example, GPT-3.5 with 3 samples could beat GPT-4 with a single sample. Given the $30\times$ cost differences ⁴, it is worth considering using GPT-3.5 with multiple samples as an alternative to GPT-4 with a single sample. Similarly, PaLM-2 could reach the GPT-3.5 level after 7 samples, again confirming the effectiveness of this method.

4.3 Few-shot Prompting

Recent studies have shown the great potential of few-shot prompting across diverse tasks including text classification (Wang et al., 2021; Zhang et al., 2022; Min et al., 2021), logical reasoning (Han et al., 2022), code generation (Bareiß et al., 2022; Joshi et al., 2023). To evaluate the effectiveness of few-shot prompting in cloud configuration generation, we provide the 3 example question-answer pairs as shown in Appendix C, each as one shot, in the prompt and evaluate the unit test scores over the original dataset. The results are shown in Table 6.

Surprisingly, our findings suggest that few-shot prompting does not yield significant improvements in unit test scores across varying model sizes. The reason could relate to the

⁴As of Oct. 1, 2023, the cost for GPT-3.5 turbo with 4k context is \$0.002 per 1k output tokens, while the cost for GPT-4 with 8k context is \$0.06 per 1k tokens.



(a) Predicted scores of 12 models on the unit test classifier. (b) SHAP values of different features in the prediction.

Figure 9: Analysis of the unit test classifier using text-level and YAML-aware scores.

heterogeneity of the cloud-native knowledge domain. For example, Kubernetes configurations are distinguished by their `kind` fields, each having unique formats and essential fields. Given this diversity, few-shot prompting may cover only a limited subset of configuration types, leaving a vast majority unaddressed. Further research is needed to reveal the underlying factors and refine the methodology for cloud-native code generation settings.

4.4 Predicting Unit Test Results

As we have demonstrated in §3.3 and §3.4, running the unit tests comes with both time and financial costs. Can we spare the costs by predicting unit test scores from text-level and YAML-aware scores?

To answer the question, we collect over 4000 YAML files and their evaluation scores, generated by 12 distinct models in Table 5. Using these scores, we train an XGBoost-based classifier (Chen & Guestrin, 2016) to predict whether a specific YAML file would pass the unit test, considering both text-level and YAML-aware scores as input features. We simulate scenarios in which we evaluate new models without running unit tests by masking the target unit scores and training the classifier on the data from the rest 11 models. The result is shown in Figure 9(a). The relative order could be maintained for most cases by the classifier except for `wizard-34b` and `codellama-7b`. However, the relative error could be as high as 80%, with most errors between 5% to 30%. In conclusion, the classifier provides a rough understanding of the ranking, but we still need to run unit tests for more accurate evaluation.

We also perform SHAP (SHapley Additive exPlanations) (Lundberg et al., 2020) to further analyze the significance of each feature. The result is shown in Figure 9(b), indicating that key-value wildcard match is the most important and accurate feature, while BLEU and Edit distance are important but noisy metrics.

Table 7: Comparison of CloudEval-YAML to other benchmarks for code generation.

Dataset	Problem Domain	Special Eval. Metric ¹	# of Problems	Data Source	Natural Lang.
HumanEval	Python algorithm	Unit tests	164	Hand-written	EN
MBPP	Basic Python	Unit tests	974	Hand-verified	EN
WikiSQL	SQL query	Execution Accuracy	88k	Hand-annotated	EN
CodeApex	C++ algorithm	Unit tests	476 ²	Online judge system	EN, ZH
MCoNaLa	Python	-	896	StackOverflow	EN, ES, JA, RU
Lyra	Python w/ embed. SQL	Code exec./AST	2000	GitHub	EN, ZH
APPS	Python	Unit tests	10k	Codeforces, Kattis	EN
CoNaLa	Python, Java	-	2879	StackOverflow	EN
Django	Python Django	Human study	19k	Django codebase	EN
Shellcode_IA32	Assembly	-	3200	shell-storm, Exploit	EN
CodeXGLUE	Python, Java	-	645k ³	Various sources	EN
CONCODE	Java classes	-	100k	GitHub repositories	EN
DS-1000	Python data science	Unit tests	1000	StackOverflow	EN
Ansible	YAML for Ansible	K-V match	112k	GitHub, GitLab	EN
CloudEval-YAML	YAML for Cloud apps	Unit tests, K-V wildcard	1011	Hand-written (337/1011)	EN, ZH

¹ We exclude widely used text-level evaluation metrics such as exact match and BLEU.

² We include the problems of the code generation task only, excluding the programming comprehension task.

³ We include the text-code category only, excluding code-code, code-text and text-text.

5 RELATED WORKS

We compare CloudEval-YAML to other benchmarks as shown in Table 7. CloudEval-YAML is a hand-written dataset that allows us to customize to specific domains and focus on real-world problems, leading to a more realistic evaluation. There are several other hand-written benchmarks, including HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021) and WikiSQL (Zhong et al., 2017). HumanEval contains 164 hand-written Python programming problems that benchmark language comprehension, algorithms, and simple mathematics, and MBPP (Austin et al., 2021) contains 974 entry-level Python problems, while WikiSQL (Zhong et al., 2017) is a hand-annotated benchmark dataset aimed at converting natural language queries into SQL queries. None of these benchmarks is targeting cloud applications.

Aside from that, there are many non-hand written datasets derived from online sources, including APPS (Hendrycks et al., 2021), CoNaLa (Yin et al., 2018; Orlanski & Gittens, 2021), Django (Oda et al., 2015), BinSum (Jin et al., 2023), Shellcode_IA32 (Liguori et al., 2021), CodeXGLUE (Lu et al., 2021), DRB-ML (Chen et al., 2023), HPC-GPT (Ding et al., 2023), CONCODE (Iyer et al., 2018), DS-1000 (Lai et al., 2023), and Ansible-YAML (Pujar et al., 2023). Among them, DS-1000 (Lai et al., 2023) is a dataset of Python data science problems collected and adapted from StackOverflow. Ansible-YAML (Pujar et al., 2023) focuses on the development of Ansible Wisdom, a natural language to Ansible-YAML code generation tool, both works inspire the design of CloudEval-YAML, but we choose to focus on the hand-written, realistic problems that focus on real

cloud application and are not common in the public datasets.

Some benchmark datasets have been created to support multiple natural languages. CodeApex (Fu et al., 2023) assesses the performance of LLMs involving bilingual prompt strategies. MCoNaLa (Wang et al., 2022) is another multilingual dataset modeled off of the methodology from CoNaLa (Yin et al., 2018) while extending beyond English. Lyra (Liang et al., 2021) is a dataset in Python with embedded SQL for code generation tasks. Each program comes with both Chinese and English comments.

6 CONCLUSION

We present the CloudEval-YAML benchmark, consisting of the first hand-written dataset that includes more than 1000 realistic problems for YAML configuration generation in cloud applications, accompanied by a complete end-to-end evaluation platform and functional correctness evaluation. We present the design of a scalable, automated evaluation platform featuring a VM cluster that can evaluate the generated code efficiently across various metrics that lead to various observations to facilitate selecting, evaluating, and optimizing code generation models.

In the future, we plan to tackle longer, more complex YAML files with AI agents that break the goal into sub-tasks such as AutoGPT (aut, 2023). We also intend to extend CloudEval-YAML to more cloud applications and programming languages. Additionally, we aim to experiment with state-of-the-art prompting methods such as chain-of-thought (Wei et al., 2022), retrieval-based LLM (Lewis et al., 2020), and fine-tuning techniques (Gu et al., 2021).

7 ACKNOWLEDGEMENT

Yifei Xu, Yuning Chen and Xumiao Zhang were supported by Alibaba Group through Alibaba Research Intern Program. We also sincerely thank Hang Qiu, Yang Li, Jian Tan, and all anonymous reviewers for their valuable feedback.

REFERENCES

- Autogpt: the heart of the open-source agent ecosystem. <https://github.com/Significant-Gravitas/AutoGPT>, 2023.
- ChatGPT. <https://openai.com/chatgpt>, 2023.
- cloud computing market size. <https://www.marketsandmarkets.com/Market-Reports/cloud-computing-market-234.html>, 2023.
- Cncf cloud native survey interactive landscape. <https://landscape.cncf.io/>, 2023.
- Codex. <https://openai.com/blog/openai-codex>, 2023.
- Copilot. <https://github.com/features/copilot>, 2023.
- GPT-4. <https://openai.com/gpt-4>, 2023.
- Welcome! | minikube. <https://minikube.sigs.k8s.io/docs/>, 2023.
- PaLM 2. <https://ai.google/discover/palm2>, 2023.
- Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Bareiß, P., Souza, B., d’Amorim, M., and Pradel, M. Code generation tools (almost) for free? a study of few-shot, pre-trained language models on code. *arXiv preprint arXiv:2206.01335*, 2022.
- Bird, S., Klein, E., and Loper, E. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O’Reilly Media, Inc.", 2009.
- Chen, L., Ding, X., Emani, M., Vanderbruggen, T., Lin, P.-h., and Liao, C. Data race detection using large language models. *arXiv preprint arXiv:2308.07505*, 2023.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Ding, B., Qin, C., Zhao, R., Luo, T., Li, X., Chen, G., Xia, W., Hu, J., Luu, A. T., and Joty, S. Data augmentation using llms: Data perspectives, learning paradigms and challenges. *arXiv preprint arXiv:2403.02990*, 2024.
- Ding, X., Chen, L., Emani, M., Liao, C., Lin, P.-H., Vanderbruggen, T., Xie, Z., Cerpa, A. E., and Du, W. Hpc-gpt: Integrating large language model for high-performance computing. In *Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis (SC-W 2023)*, 2023.
- Feng, Y. and Jeon, H. Understanding scalability of multi-gpu systems. In *Proceedings of the 15th Workshop on General Purpose Processing Using GPU*, pp. 36–37, 2023.
- Feng, Y., Jeon, H., Blagojevic, F., Guyot, C., Li, Q., and Li, D. Attmemo: Accelerating transformers with memoization on big memory systems. *arXiv preprint arXiv:2301.09262*, 2023.
- Fu, L., Chai, H., Luo, S., Du, K., Zhang, W., Fan, L., Lei, J., Rui, R., Lin, J., Fang, Y., et al. Codeapex: A bilingual programming evaluation benchmark for large language models. *arXiv preprint arXiv:2309.01940*, 2023.
- Gu, Y., Han, X., Liu, Z., and Huang, M. Ppt: Pre-trained prompt tuning for few-shot learning. *arXiv preprint arXiv:2109.04332*, 2021.
- Han, S., Schoelkopf, H., Zhao, Y., Qi, Z., Riddell, M., Benson, L., Sun, L., Zubova, E., Qiao, Y., Burtell, M., et al. Folio: Natural language reasoning with first-order logic. *arXiv preprint arXiv:2209.00840*, 2022.
- Hendrycks, D., Basart, S., Kadavath, S., Mazeika, M., Arora, A., Guo, E., Burns, C., Puranik, S., He, H., Song, D., et al. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.
- Iyer, S., Konstas, I., Cheung, A., and Zettlemoyer, L. Mapping language to code in programmatic context. *arXiv preprint arXiv:1808.09588*, 2018.

- Jin, X. and Wang, Y. Understand legal documents with contextualized large language models. *arXiv preprint arXiv:2303.12135*, 2023.
- Jin, X., Larson, J., Yang, W., and Lin, Z. Binary code summarization: Benchmarking chatgpt/gpt-4 and other large language models. *arXiv preprint arXiv:2312.09601*, 2023.
- Joshi, H., Sanchez, J. C., Gulwani, S., Le, V., Verbruggen, G., and Radiček, I. Repair is nearly generation: Multilingual program repair with llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 5131–5140, 2023.
- Kabir, S., Udo-Imeh, D. N., Kou, B., and Zhang, T. Who answers it better? an in-depth analysis of chatgpt and stack overflow answers to software engineering questions. *arXiv preprint arXiv:2308.02312*, 2023.
- Kulal, S., Pasupat, P., Chandra, K., Lee, M., Padon, O., Aiken, A., and Liang, P. S. Spoc: Search-based pseudocode to code. *Advances in Neural Information Processing Systems*, 32, 2019.
- Lai, Y., Li, C., Wang, Y., Zhang, T., Zhong, R., Zettlemoyer, L., Yih, W.-t., Fried, D., Wang, S., and Yu, T. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pp. 18319–18345. PMLR, 2023.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- Liang, Q., Sun, Z., Zhu, Q., Zhang, W., Yu, L., Xiong, Y., and Zhang, L. Lyra: A benchmark for turducken-style code generation. *arXiv preprint arXiv:2108.12144*, 2021.
- Liguori, P., Al-Hossami, E., Cotroneo, D., Natella, R., Cukic, B., and Shaikh, S. Shellcode_ia32: A dataset for automatic shellcode generation. *arXiv preprint arXiv:2104.13100*, 2021.
- Lu, S., Guo, D., Ren, S., Huang, J., Svyatkovskiy, A., Blanco, A., Clement, C., Drain, D., Jiang, D., Tang, D., et al. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., and Lee, S.-I. From local explanations to global understanding with explainable ai for trees. *Nature machine intelligence*, 2(1):56–67, 2020.
- Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., and Jiang, D. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- Min, S., Lewis, M., Hajishirzi, H., and Zettlemoyer, L. Noisy channel language model prompting for few-shot text classification. *arXiv preprint arXiv:2108.04106*, 2021.
- Moritz, P., Nishihara, R., Wang, S., Tumanov, A., Liaw, R., Liang, E., Elibol, M., Yang, Z., Paul, W., Jordan, M. I., et al. Ray: A distributed framework for emerging {AI} applications. In *13th USENIX symposium on operating systems design and implementation (OSDI 18)*, pp. 561–577, 2018.
- Oda, Y., Fudaba, H., Neubig, G., Hata, H., Sakti, S., Toda, T., and Nakamura, S. Learning to generate pseudo-code from source code using statistical machine translation. In *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 574–584. IEEE, 2015.
- Orlanski, G. and Gittens, A. Reading stackoverflow encourages cheating: adding question text improves extractive code generation. *arXiv preprint arXiv:2106.04447*, 2021.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. 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, 2002.
- Pujar, S., Buratti, L., Guo, X., Dupuis, N., Lewis, B., Suneja, S., Sood, A., Nalawade, G., Jones, M., Morari, A., et al. Automated code generation for information technology tasks in yaml through large language models. *arXiv preprint arXiv:2305.02783*, 2023.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X. E., Adi, Y., Liu, J., Remez, T., Rapin, J., et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

- Vaithilingam, P., Zhang, T., and Glassman, E. L. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In *Chi conference on human factors in computing systems extended abstracts*, pp. 1–7, 2022.
- Wang, C., Wang, J., Qiu, M., Huang, J., and Gao, M. Transprompt: Towards an automatic transferable prompting framework for few-shot text classification. In *Proceedings of the 2021 conference on empirical methods in natural language processing*, pp. 2792–2802, 2021.
- Wang, Z., Cuenca, G., Zhou, S., Xu, F. F., and Neubig, G. Mconala: a benchmark for code generation from multiple natural languages. *arXiv preprint arXiv:2203.08388*, 2022.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022.
- Whitehouse, C., Choudhury, M., and Aji, A. F. Llm-powered data augmentation for enhanced crosslingual performance. *arXiv preprint arXiv:2305.14288*, 2023.
- Xu, Y., Chen, Y., Zhang, X., Lin, X., Hu, P., Ma, Y., Lu, S., Du, W., Mao, Z. M., Zhai, E., et al. Cloudeval-yaml: A realistic and scalable benchmark for cloud configuration generation. 2023.
- Yang, K., Chen, Y., and Du, W. OrchLoc: In-Orchard Localization via a Single LoRa Gateway and Generative Diffusion Model-based Fingerprinting. In *ACM MobiSys*, 2024.
- Yin, P., Deng, B., Chen, E., Vasilescu, B., and Neubig, G. Learning to mine aligned code and natural language pairs from stack overflow. In *Proceedings of the 15th international conference on mining software repositories*, pp. 476–486, 2018.
- Zhang, H., Zhang, X., Huang, H., and Yu, L. Prompt-based meta-learning for few-shot text classification. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 1342–1357, 2022.
- Zhong, V., Xiong, C., and Socher, R. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv preprint arXiv:1709.00103*, 2017.

Table 8: Statistics of YAML files in top-100 most starred cloud native applications.

Repo Name	Github Stars	Total Files	YAML Files	Repo Name	Github Stars	Total Files	YAML Files	Repo Name	Github Stars	Total Files	YAML Files
GitLab	23368	58372	4721	Dgraph	19620	2231	161	Terraform	38875	5704	36
Kubernetes	101881	29662	4715	Salt Project	13513	7242	153	Flink	21993	27228	30
Elastic	65213	35747	3143	Docker Compose	30543	466	147	Apollo	28360	1512	28
GraphQL	30135	13667	2169	Vitess	16897	5579	142	gVisor	14172	3723	26
Istio	33694	6261	2081	containerd	14857	6523	138	Sentinel	21422	3487	25
Ansible	58659	7236	1914	Serverless	45187	1805	131	go-zero	25550	1382	22
ShardingSphere	18807	21945	1632	CockroachDB	27828	18499	118	Seata	24226	3904	21
llvm	21975	148442	1202	k3s	24517	750	97	Packer	14612	1450	20
Argo	14145	4172	1118	Logstash	13639	3835	88	Wasmer	16300	2007	19
Skaffold	14219	16345	1044	Apache Spark	36800	24415	85	Portainer	26644	3063	19
Kubespray	14472	2093	900	Kong	35947	1888	75	Golang	114620	14022	18
SkyWalking	22442	5999	802	SST	17715	4683	73	SOPS	13823	190	18
Cilium	16516	19972	780	Rust	85579	46998	69	Redis	61572	1679	16
MongoDB	24425	49784	743	gRPC	39066	12629	68	kratos	21387	861	16
Backstage	23285	12300	613	Vault	27546	9175	66	NATS	24451	580	16
Grafana Loki	20163	15520	554	DragonflyDB	21064	615	64	Zig	26009	16173	15
Helm	24953	1784	540	Consul	26921	13084	62	Jenkins	21453	13139	15
Envoy	22759	13470	520	Keycloak	17472	14535	59	Apache Hadoop	13858	9562	14
Pulumi	17622	8179	467	Presto	15087	13493	57	Dubbo	39400	5399	14
Teleport	14225	8884	419	InfluxData	26133	2007	56	TiDB	34880	6235	14
Traefik	44719	1870	339	ORY Hydra	14434	2556	56	OpenFaaS	23512	1100	14
minikube	27261	2368	316	OpenAPI	27136	181	55	emscripten	24266	9596	11
SlimToolkit	17269	6545	305	Sentry	35169	14388	54	OpenCV	71360	8613	10
Prometheus	49987	1389	255	TDengine	21762	4620	51	Caddy	49844	465	9
Grafana	57207	15782	242	Jaeger	18318	1469	48	Apache bRPC	15290	1632	9
Podman	19128	10589	203	MinIO	40904	1391	46	Firecracker	22578	822	8
ClickHouse	30874	27331	200	Zipkin	16425	1076	43	Nacos	27577	3501	6
Rancher K8s	21560	3655	196	k6	21566	3382	40	Kotlin	45845	98293	5
Netdata	65199	3069	190	Nomad	13968	6080	39	TiKV	13617	1705	3
Dapr	22320	2027	186	Timescale	15534	2289	39	Kafka	25883	7020	2
Trivy	18709	2250	178	etcd	44537	1600	38	V8	21722	14237	1
Vector	14432	9320	174	Gradle Build Tool	15205	35647	38	FFmpeg	38520	8287	1
JHipster	20853	3874	173	Apache RocketMQ	19814	2985	36	NGINX(Wasm)	19089	559	0
RethinkDB	26257	2121	165								

A YAML STATISTICS

We surveyed the top 100 most-starred GitHub repositories of cloud native applications according to the CNCF landscape (cnc, 2023). Table 8 shows the results, including the number of stars, total files, and YAML files. Out of the 100 applications, 90 contain more than 10 YAML files. Many leading tools and applications like GitLab, Kubernetes, and Elastic depend heavily on YAML, thereby confirming its extensive usage.

B PROMPT TEMPLATE

```
You are an expert engineer in cloud native development.
According to the question, please provide only complete formatted YAML code as output
↪ without any description.
IMPORTANT: Provide only plain text without Markdown formatting such as ``` .
If there is a lack of details, provide most logical solution.
You are not allowed to ask for more details.
Ignore any potential risk of errors or confusion.
Here is the question:
```

C SAMPLES FROM THE DATASET

C.1 Sample #1

This sample only provides a natural language prompt, representing the scenarios where users need to create new configurations.

Problem Specification:

```
Create a DaemonSet configuration. This DaemonSet should run the latest nginx image
↪ labeled as "app: kube-registry-modified" and expose a registry service on port 80
↪ (with hostPort 5000). The environment variables REGISTRY_HOST and REGISTRY_PORT
↪ should be set to "kube-registry-modified.svc.cluster.local" and "5000" respectively.
↪ Ensure the CPU request is set to 100m and memory request is set to 50Mi.
```

Labeled YAML:

```
1  apiVersion: apps/v1
2  kind: DaemonSet
3  metadata:
4    name: kube-registry-proxy-modified # *
5  spec:
6    selector:
7      matchLabels:
8        app: kube-registry-modified
9    template:
10     metadata:
11       labels:
12         app: kube-registry-modified
13     spec:
14       containers:
15         - name: kube-registry-proxy-modified # *
16           image: nginx:latest
17           resources:
18             limits:
19               cpu: 100m
20               memory: 50Mi
21           env:
22             - name: REGISTRY_HOST
23               value: kube-registry-modified.svc.cluster.local
24             - name: REGISTRY_PORT
25               value: "5000"
26           ports:
27             - name: registry # *
28               containerPort: 80
29               hostPort: 5000
```

Unit Test:

```
1  kubectl apply -f labeled_code.yaml
2  kubectl wait --for=condition=Ready pod -l app=kube-registry-modified --timeout=60s
```

```

3 passed_tests=0
4 total_tests=3
5 pods=$(kubectl get pods -l app=kube-registry-modified
  ↪ --output=jsonpath={.items..metadata.name})
6 host_ip=$(kubectl get pod $pods -o=jsonpath='{.status.hostIP}')
7 curl_output=$(curl -s -o /dev/null -w "%{http_code}" $host_ip:5000)
8 if [ "$curl_output" == "200" ]; then
9     ((passed_tests++))
10 else
11     exit 1
12 fi
13 env_vars=$(kubectl get pods --selector=app=kube-registry-modified
  ↪ -o=jsonpath='{.items[0].spec.containers[0].env[*].name}')
14 if [[ $env_vars == *"REGISTRY_HOST"* && $env_vars == *"REGISTRY_PORT"* ]]; then
15     ((passed_tests++))
16 fi
17 cpu_limit=$(kubectl get pod $pods
  ↪ -o=jsonpath='{.spec.containers[0].resources.limits.cpu}')
18 memory_limit=$(kubectl get pod $pods
  ↪ -o=jsonpath='{.spec.containers[0].resources.limits.memory}')
19 if [ "$cpu_limit" == "100m" ] && [ "$memory_limit" == "50Mi" ]; then
20     ((passed_tests++))
21 fi
22 if [ $passed_tests -eq $total_tests ]; then
23     echo unit_test_passed
24 fi

```

C.2 Sample #2

This sample provides a context YAML and seeks functionality extensions.

Problem Specification:

Given the following YAML, please help me create a service with load balancer that uses
 ↪ the nginx selector, exposed on port 80.
 It should be accessible via browser.

```

1 apiVersion: apps/v1
2 kind: Deployment
3 metadata:
4     name: nginx-deployment
5 spec:
6     replicas: 3
7     selector:
8         matchLabels:
9             app: nginx
10    template:
11        metadata:
12            labels:
13                app: nginx
14        spec:
15            containers:
16                - name: nginx-container

```



```
17     image: nginx:latest
18     ports:
19     - containerPort: 80
```

Labeled YAML:

```
1  apiVersion: v1
2  kind: Service
3  metadata:
4    name: nginx-service    # *
5  spec:
6    selector:
7      app: nginx
8    ports:
9      - name: http
10       port: 80
11       targetPort: 80
12     type: LoadBalancer
```

Unit Test:

```
1  echo "apiVersion: apps/v1
2  kind: Deployment
3  [... the same as the YAML context in problem specification, omitted by brevity]
4    - containerPort: 80" | kubectl apply -f -
5  kubectl wait --for=condition=ready deployment --all --timeout=15s
6  kubectl apply -f labeled_code.yaml
7  sleep 15
8  kubectl get svc
9  timeout -s INT 8s minikube service nginx-service > bash_output.txt 2>&1
10 cat bash_output.txt
11 grep "Opening service default/nginx-service in default browser..." bash_output.txt &&
   ↪ echo unit_test_passed
```

C.3 Sample #3

This is a debugging query sourced from StackOverflow. The raw YAML is included along with the error report. Typically, the response includes an error analysis unless restricted by the prompt template.

Problem Specification:

Given the following YAML which is not functionally correct:

```
1  apiVersion: networking.k8s.io/v1
2  kind: Ingress
3  metadata:
4    name: test-ingress
5  annotations:
6    nginx.ingress.kubernetes.io/rewrite-target: /
```

```
7 spec:
8   rules:
9     - http:
10       paths:
11         - path: /
12           backend:
13             serviceName: test-app
14             servicePort: 5000
```

When executing it, it would report the error:

```
Error from server (BadRequest): error when creating "wrong.yaml": Ingress in version "v1"
↳ cannot be handled as a Ingress: strict decoding error: unknown field "annotations",
↳ unknown field "spec.rules[0].http.paths[0].backend.serviceName", unknown field
↳ "spec.rules[0].http.paths[0].backend.servicePort"
```

Please debug it to make it valid. Please provide the entire YAML.

Labeled YAML:

```
1 apiVersion: networking.k8s.io/v1
2 kind: Ingress
3 metadata:
4   name: minimal-ingress
5   annotations:
6     nginx.ingress.kubernetes.io/rewrite-target: /
7 spec:
8   rules:
9     - http:
10       paths:
11         - path: /
12           pathType: Prefix
13           backend:
14             service:
15               name: test-app
16               port:
17                 number: 5000
```

Unit Test:

```
1 kubectl apply -f labeled_code.yaml
2 kubectl wait --namespace default --for=condition=SYNCD ingress --all --timeout=15s
3 kubectl describe ingress minimal-ingress | grep "test-app:5000" && echo unit_test_passed
```

D PROMPTS FOR DATA AUGMENTATION

D.1 Prompt for question simplification

You are a helpful assistant trying to rephrase INPUT paragraph that includes a question.
 ↪ Don't answer the question directly. Instead, please rewrite the question to make it
 ↪ as concise as possible, using abbreviations but include necessary informations. The
 ↪ question does not have to be a complete sentence. Don't modify code surrounded by ```.
 ↪ Let me first provide a few examples.

Example #1:

INPUT: Craft a yaml file to define a Kubernetes LimitRange. Containers within the cluster
 ↪ should have a default CPU request of 100m and a memory request of 200Mi. Any Pod
 ↪ created should not exceed a maximum CPU usage of 150m or a memory usage of 250Mi.
 OUTPUT: Create a yaml for a k8s LimitRange. Containers CPU/Mem default to 100m/200Mi,
 ↪ pods max to 150m/250Mi.

Example #2:

INPUT: Please write a YAML file that defines firstly a Service and then a Deployment. The
 ↪ Deployment runs a single MySQL instance using the latest image on port 3306, with the
 ↪ environment MYSQL_ROOT_PASSWORD=password. The deployment should also define a volume
 ↪ mount for /var/lib/mysql, and then creates a PersistentVolumeClaim that matches the
 ↪ existing one. The Service simply expose the deployment on its port. All potential
 ↪ names should be mysql and labels should be app: mysql.
 OUTPUT: Write a YAML defining a Service & Deployment. Deployment runs a MySQL instance on
 ↪ port 3306, env MYSQL_ROOT_PASSWORD=password, volume mount /var/lib/mysql, creates a
 ↪ matching PersistentVolumeClaim. Service exposes the deployment on its port. Using
 ↪ names mysql & labels app: mysql.

Example #3:

INPUT:
 Is there a way to provide environment variables from a file in Kubernetes, for example
 ↪ when defining a pod, instead of hardcoding them like this:

```

1  apiVersion: v1
2  kind: Pod
3  metadata:
4    labels:
5      context: docker-k8s-lab
6      name: mysql-pod
7  name: mysql-pod
8  spec:
9    containers:
10   - env:
11     - name: MYSQL_USER
12       value: mysql
13     - name: MYSQL_PASSWORD
14       value: mysql
15     - name: MYSQL_DATABASE
16       value: sample
17     - name: MYSQL_ROOT_PASSWORD
18       value: supersecret
19   image: "mysql:latest"
20   name: mysql

```

```
21     ports:
22       - containerPort: 3306
```

Please provide the entire YAML, supposing there is a Secret named mysql-secret, with all ↪ values in it.

OUTPUT:

Can k8s use env var from a file instead of hardcoding? Given:

```
1  apiVersion: v1
2  kind: Pod
3  metadata:
4    labels:
5      context: docker-k8s-lab
6      name: mysql-pod
7  name: mysql-pod
8  spec:
9    containers:
10     - env:
11       - name: MYSQL_USER
12         value: mysql
13       - name: MYSQL_PASSWORD
14         value: mysql
15       - name: MYSQL_DATABASE
16         value: sample
17       - name: MYSQL_ROOT_PASSWORD
18         value: supersecret
19     image: "mysql:latest"
20     name: mysql
21     ports:
22       - containerPort: 3306
```

Assume a Secret named mysql-secret with all values. Provide the full YAML.

Example #4:

INPUT:

I'm working with the bookinfo application in our Istio setup.

I recall there was a DestinationRule specifically for the ratings service in the prod ↪ namespace, which ensures traffic is load balanced using the LEAST_REQUEST strategy. Please provide me the exact configuration for that.

OUTPUT:

Provide exact config for DestinationRule for ratings service in prod namespace using ↪ LEAST_REQUEST lb strategy in Istio/bookinfo app setup.

Example #5:

INPUT:

I need a Istio destination rule YAML set up for the bookinfo application's ratings ↪ service in the prod namespace.

This rule had the main traffic load balanced using the LEAST_REQUEST strategy.

Additionally, there was a specific subset named testversion using version v3 labels, and ↪ for this subset, the traffic was load balanced with a ROUND_ROBIN approach.

Please provide me the entire YAML configuration for this.

OUTPUT:

Provide Istio DestinationRule YAML for bookinfo app's ratings service in prod ns. Main ↪ traffic uses LEAST_REQUEST lb, subset "testversion" uses labels v3 and ROUND_ROBIN lb ↪ strategy.

That's all the examples I have. Now let's working on the following input:
INPUT: {question}

D.2 Prompt for question translation

You are a helpful assistant trying to translate INPUT paragraph that includes a question
↪ into fluent Chinese in a developer's tone.

Don't answer the question directly. Instead, please just translate the question for a
↪ Chinese developer to answer. The question does not have to be a formal and complete
↪ sentence. Don't modify code surrounded by ```. Let me first provide a few examples.

Example #1:

INPUT: Craft a yaml file to define a Kubernetes LimitRange. Containers within the cluster
↪ should have a default CPU request of 100m and a memory request of 200Mi. Any Pod
↪ created should not exceed a maximum CPU usage of 150m or a memory usage of 250Mi.

OUTPUT: 写一个 yaml 来定义 Kubernetes LimitRange。集群里的容器的默认 CPU 请求是 100m，默认内存
↪ 请求是 200Mi。创建的任何 Pod 的 CPU 使用量最大不应超过 150m，内存使用最大不应超过 250Mi。

Example #2:

INPUT: Please write a YAML file that defines firstly a Service and then a Deployment. The
↪ Deployment runs a single MySQL instance using the latest image on port 3306, with the
↪ environment MYSQL_ROOT_PASSWORD=password. The deployment should also define a volume
↪ mount for /var/lib/mysql, and then creates a PersistentVolumeClaim that matches the
↪ existing one. The Service simply expose the deployment on its port. All potential
↪ names should be mysql and labels should be app: mysql.

OUTPUT: 请写一个 YAML，先定义 Service，再定义 Deployment。该 Deployment 使用端口 3306 上的最
↪ 新映像运行单个 MySQL 实例，环境为 MYSQL_ROOT_PASSWORD=password。Deployment 还应该为
↪ /var/lib/mysql 定义一个卷挂载，然后创建一个与现有卷匹配的 PersistentVolumeClaim。该 Service
↪ 只是在其端口上公开部署。所有潜在名称定为 mysql，标签定为 app:mysql。

Example #3:

INPUT:

Is there a way to provide environment variables from a file in Kubernetes, for example
↪ when defining a pod, instead of hardcoding them like this:

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   labels:
5     context: docker-k8s-lab
6     name: mysql-pod
7 name: mysql-pod
8 spec:
9   containers:
10  - env:
11    - name: MYSQL_USER
12      value: mysql
13    - name: MYSQL_PASSWORD
14      value: mysql
15    - name: MYSQL_DATABASE
```

```
16     value: sample
17   - name: MYSQL_ROOT_PASSWORD
18     value: supersecret
19   image: "mysql:latest"
20   name: mysql
21   ports:
22     - containerPort: 3306
```

Please provide the entire YAML, supposing there is a Secret named `mysql-secret`, with all ↪ values in it.

OUTPUT:

有没有办法从 Kubernetes 中的文件提供环境变量，比如在定义 pod 时，而不是像下面这样对它们进行硬编码：

```
1  apiVersion: v1
2  kind: Pod
3  metadata:
4    labels:
5      context: docker-k8s-lab
6      name: mysql-pod
7  name: mysql-pod
8  spec:
9    containers:
10     - env:
11       - name: MYSQL_USER
12         value: mysql
13       - name: MYSQL_PASSWORD
14         value: mysql
15       - name: MYSQL_DATABASE
16         value: sample
17       - name: MYSQL_ROOT_PASSWORD
18         value: supersecret
19     image: "mysql:latest"
20     name: mysql
21     ports:
22     - containerPort: 3306
```

请提供整个 YAML，假设有一个名为 `mysql-secret` 的 Secret，其中包含所有值。

Example #4:

INPUT:

I'm working with the bookinfo application in our Istio setup.

I recall there was a DestinationRule specifically for the ratings service in the prod ↪ namespace, which ensures traffic is load balanced using the LEAST_REQUEST strategy.

Please provide me the exact configuration for that.

OUTPUT:

我正在 Istio 配置中使用 bookinfo 应用。

我记得有一个专门用于生产命名空间中的 ratings 服务的 DestinationRule，它确保使用 LEAST_REQUEST 策略 ↪ 进行流量负载均衡。

请为此提供确切的配置。

Example #5:

INPUT:

I need a Istio destination rule YAML set up for the bookinfo application's ratings ↪ service in the prod namespace.

This rule had the main traffic load balanced using the LEAST_REQUEST strategy.

Additionally, there was a specific subset named testversion using version v3 labels, and ↪ for this subset, the traffic was load balanced with a ROUND_ROBIN approach. Please provide me the entire YAML configuration for this.

OUTPUT:

我需要为 prod namespace 中的 bookinfo 应用的评级服务配置一个 Istio destination rule YAML。该规则使用 LEAST_REQUEST 策略对主要流量进行负载均衡。
此外, 还有一个名为 testversion 的特定子集, 使用版本 v3 标签, 对于该子集, 流量通过 ROUND_ROBIN 方法进行负载均衡。
请为此提供完整的 YAML 配置。

That's all the examples I have. Now let's working on the following input:
INPUT: {question}

The data augmentation part is an enhancement to the early version of CloudEval-YAML (Xu et al., 2023). Here GPT-4 only acts as an assistant while we still keep human experts as the golden standard, and employ heavy manual revision in this process. From Table 5, we can see that augmentation does not bring significant benefits for GPT-4 itself: GPT-3.5 closes the gap by 16 in the simplified dataset; LLaMa-70B also closes the gap by 3 in the translated dataset. It aligns with common practices for generative data augmentation (Whitehouse et al., 2023; Ding et al., 2024; Yang et al., 2024; Jin & Wang, 2023).

E PERFORMANCE ANALYSIS ON DIFFERENT AFFECTING FACTORS

Table 9 shows the detailed model performance on several key factors. The visualizations are in Figure 6.

Table 9: Performance analysis on different affecting factors

	Category			Code context		Reference code length			Question token number		
	Kubernetes	Envoy	Istio	w/ code	w/o code	[0, 15)	[15, 30)	≥30	[0, 50)	[50, 100)	≥100
gpt-4	0.601	0.1	0.385	0.51	0.547	0.625	0.616	0.237	0.619	0.563	0.427
gpt-3.5	0.466	0.122	0.385	0.429	0.416	0.534	0.477	0.169	0.476	0.423	0.378
PaLM-2-bison	0.406	0.05	0.231	0.354	0.358	0.455	0.413	0.118	0.44	0.345	0.309
Llama-2-70b-chat	0.099	0.049	0	0.082	0.095	0.216	0.058	0.013	0.167	0.077	0.045
Llama-2-13b-chat	0.085	0.049	0	0.061	0.089	0.125	0.081	0.013	0.143	0.063	0.045
Wizardcoder-34b-v1.0	0.067	0.05	0.231	0.068	0.074	0.159	0.052	0.013	0.131	0.049	0.055
Llama-2-7b-chat	0.039	0.05	0	0.054	0.026	0.08	0.029	0.013	0.036	0.042	0.036
Wizardcoder-15b-v1.0	0.032	0.049	0.077	0.061	0.016	0.045	0.041	0.013	0.012	0.035	0.054
Llama-7b	0.035	0.05	0	0.054	0.021	0.057	0.035	0.013	0.024	0.042	0.036
Llama-13b-lora	0.021	0.049	0	0.041	0.011	0.034	0.017	0.026	0.012	0.035	0.054
Codellama-13b-instruct	0.011	0.05	0	0.007	0.021	0.034	0.006	0.013	0	0.021	0.018
Codellama-7b-instruct	0.007	0.049	0.077	0.014	0.016	0.034	0.006	0.013	0.012	0.007	0.027