

INSEVA: A COMPREHENSIVE CHINESE BENCHMARK FOR LARGE LANGUAGE MODELS IN INSURANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

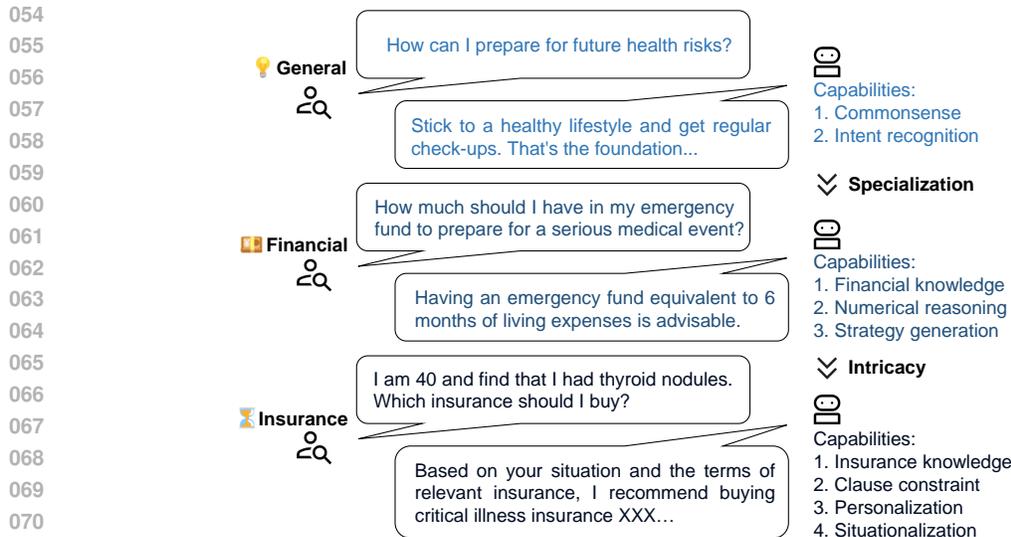
Insurance, as a critical component of the global financial system, demands high standards of accuracy and reliability in AI applications. While existing benchmarks evaluate AI capabilities across various domains, they often fail to capture the unique characteristics and requirements of the insurance domain. To address this gap, we present INSEva, a comprehensive Chinese benchmark specifically designed for evaluating AI systems' knowledge and capabilities in insurance. INSEva features a multi-dimensional evaluation taxonomy covering business areas, task formats, difficulty levels, and cognitive-knowledge dimension, comprising 38,704 high-quality evaluation examples sourced from authoritative materials. Our benchmark implements tailored evaluation methods for assessing both faithfulness and completeness in open-ended responses. Through extensive evaluation of 9 state-of-the-art Large Language Models (LLMs), we identify significant performance variations across different dimensions. While general LLMs demonstrate basic insurance domain competency with average scores above 80, substantial gaps remain in handling complex, real-world insurance scenarios. The benchmark will be public soon.

1 INTRODUCTION

Insurance is one of the fundamental pillars of the global financial ecosystem (Pfeifer & Langen, 2021), performing irreplaceable infrastructure functions in economic stability worldwide. Recently, Artificial Intelligence (AI) has emerged as a transformative force in this field (Jaiswal, 2023). Due to the high risk of insurance decisions, AI systems are required not only to possess domain knowledge but also to ensure high accuracy and harmlessness. Therefore, specialized benchmarks for evaluating the domain-specific knowledge and capabilities of AI in insurance have become particularly crucial.

Many recent studies have attempted to construct benchmarks in the financial domain, which is the parent domain of insurance, yet few focus specifically on insurance. Existing financial benchmarks (Chen et al., 2021; Zhu et al., 2021; Chen et al., 2022; Shah et al., 2022; Zhu et al., 2024; Zhang et al., 2023; Peng et al., 2025; Xie et al., 2024) evaluate LLMs' capabilities within the broad financial domain. The financial domain centers on capital circulation, value appreciation, and financial instrument transactions, emphasizing wealth growth through capital operations. Insurance, however, focuses on risk management, risk assessment, pricing, and claims services under complex constraints (insurance clauses) to ensure policyholder protection and economic compensation during risk events, as shown in Figure 1. These fundamental differences require a dedicated insurance benchmark, while currently only a few of the benchmarks are proposed for insurance. InsQABench (Ding et al., 2025) collects insurance commonsense question-answering data for evaluation. INSMMBench (Lin et al., 2024) focuses on evaluating cross-modal alignment capabilities in auto insurance scenarios. These benchmarks only include basic knowledge questions in insurance, which fail to adequately cover the insurance's unique risk characteristics, complex product clauses, deep causal reasoning, cross-domain knowledge integration (e.g., medical, legal), and critical elements in insurance sales processes.

To address these limitations, we present INSEva, a comprehensive Chinese benchmark specifically designed for evaluating AI systems' knowledge and capabilities in the insurance domain. Based on an in-depth study of the insurance production environment, we develop a multi-dimensional evaluation taxonomy that encompasses business areas, task formats, difficulty levels, and cognitive-



072 Figure 1: A comparison of question-answering examples and required model capabilities in the
073 general, financial, and insurance domains. It shows how questions are becoming increasingly spe-
074 cialized and intricate in the general, financial, and insurance domains, and how this places different
075 demands on LLM capabilities.

076
077
078
079 knowledge domains, thereby ensuring comprehensive assessment coverage. Our data construction
080 pipeline incorporates three key stages: (1) collection from authoritative sources, (2) systematic data
081 augmentation to enhance diversity while maintaining domain authenticity, and (3) rigorous quality
082 control methods involving expert validation. This methodical approach yields 38,704 high-quality
083 evaluation examples spanning the insurance domain. Given the high risk of insurance operations
084 where errors can have significant financial and regulatory consequences, we implement tailored
085 evaluation methods for different question types. For structured questions with definitive answers,
086 we employ conventional deterministic metrics (primarily accuracy), while for open-ended responses,
087 we deploy a specially designed LLM-based evaluation framework that systematically assesses both
088 faithfulness to retrieved information and completeness to ground truth. Through these methodologi-
089 cal innovations, INSEva not only provides a comprehensive assessment of insurance-specific capa-
090 bilities but also effectively differentiates model performance across multiple dimensions. The bench-
091 mark thus offers valuable guidance for domain-specific development efforts and targeted model im-
092

093 Using INSEva, we evaluate 9 state-of-the-art LLMs (including closed-source, open-source, and
094 domain-specific LLMs), revealing differentiated performance across various taxonomy dimensions
095 and identifying numerous meaningful insights. While general LLMs demonstrate a foundational
096 understanding of insurance concepts, as evidenced by average scores exceeding 80, a considerable
097 gap remains between their capabilities and the expertise required to tackle complex, real-world in-
098 surance problems. Specifically, we observe a performance bottleneck in tasks requiring logical
099 reasoning and a trade-off between faithfulness and completeness in generated content. Notably,
100 the underperformance of the financial LLMs compared to general LLMs underscores the need for
101 specialized benchmarks tailored to the unique demands of insurance, as existing financial bench-
102 marks may not adequately capture the nuances of this critical domain. Furthermore, we fine-tune an
103 LLM using data from our pipeline, achieving substantial performance gains, thereby validating the
104 effectiveness of both our data construction methodology and the benchmark itself.

105 To summarize, our contributions are as follows:

- We propose a comprehensive Chinese benchmark for the insurance domain, with a multi-
106 dimensional evaluation taxonomy and 38,704 high-quality evaluation examples from au-
107 thoritative sources.

- We design a data construction and evaluation pipeline to ensure the high-quality of data and evaluation for the insurance domain.
- We evaluate leading LLMs to show their domain-specific gaps, and further validate our benchmark’s utility by confirming it accurately reflects performance improvements from targeted fine-tuning.

2 RELATED WORKS

2.1 GENERAL BENCHMARKS

Many general benchmarks have been developed to assess the broad capabilities of models across diverse tasks and domains. For example, the GLUE benchmark (Wang et al., 2019) comprises a collection of natural language understanding tasks. MMLU (Hendrycks et al., 2021) evaluates multi-task language understanding across a wide range of subjects, assessing both knowledge acquisition and reasoning abilities. The HELM benchmark (Liang et al., 2023) adopts a holistic approach to evaluate models across a diverse set of scenarios and metrics, emphasizing real-world relevance and ethical considerations. BIG-bench (Srivastava A, 2023) introduces tasks that require reasoning, common sense, and cultural knowledge. While these general benchmarks provide valuable insights into the overall performance of LLMs, they often lack the domain-specific focus and granularity required to assess performance in specialized areas such as insurance, highlighting the need for tailored evaluation resources.

2.2 FINANCE & INSURANCE BENCHMARKS

As the parent domain of insurance, the development of financial benchmarks for financial models has accelerated in recent years, providing structured frameworks for assessing models’ capabilities in the finance domain. FinQA (Chen et al., 2021) advances numerical reasoning evaluation by constructing a dataset requiring multi-step calculations and inference from financial reports. FLUE (Shah et al., 2022) introduces comprehensive financial language understanding tasks. CFLUE (Zhu et al., 2024) evaluates models through two main dimensions in the Chinese financial domain, including knowledge assessment and application assessment. FinBen (Xie et al., 2024) is an extensive open-source evaluation benchmark, including 36 datasets spanning 24 financial tasks. Despite the proliferation of financial evaluation resources, they are difficult to apply directly to insurance, which presents unique terminology, regulatory frameworks, and reasoning requirements distinct from those in general finance contexts. Some financial benchmarks (such as CFLUE) include subtasks in insurance, but the questions in these subtasks generally focus only on basic insurance knowledge, lacking real-world applications.

Compared with the financial domain, there are far fewer benchmarks specifically for insurance. InsuranceQA (Feng et al., 2015) is the first benchmark for insurance that contains data collected from the internet. InsQABench (Ding et al., 2025) encompasses three specialized insurance QA tasks, including insurance commonsense knowledge, insurance structured database, and insurance unstructured document. Ins-MMBench (Lin et al., 2024) comprises a total of 2.2k multimodal multiple-choice questions for evaluating large vision-language models. Due to the lack of a complete insurance taxonomy, these benchmarks fail to fully cover insurance tasks.

3 INSEVA BENCHMARK

3.1 OVERVIEW

INSEva is a comprehensive resource designed to evaluate the performance of models in the insurance-specific tasks. The benchmark encompasses a diverse set of samples and tasks, drawing data from authoritative Chinese sources such as professional insurance exams, regulatory standards, and real-world business data. INSEva comprises 38,704 carefully curated examples spanning various insurance sub-domains and difficulties. Model performance is evaluated using a combination of accuracy and domain-specific metrics tailored to each task.

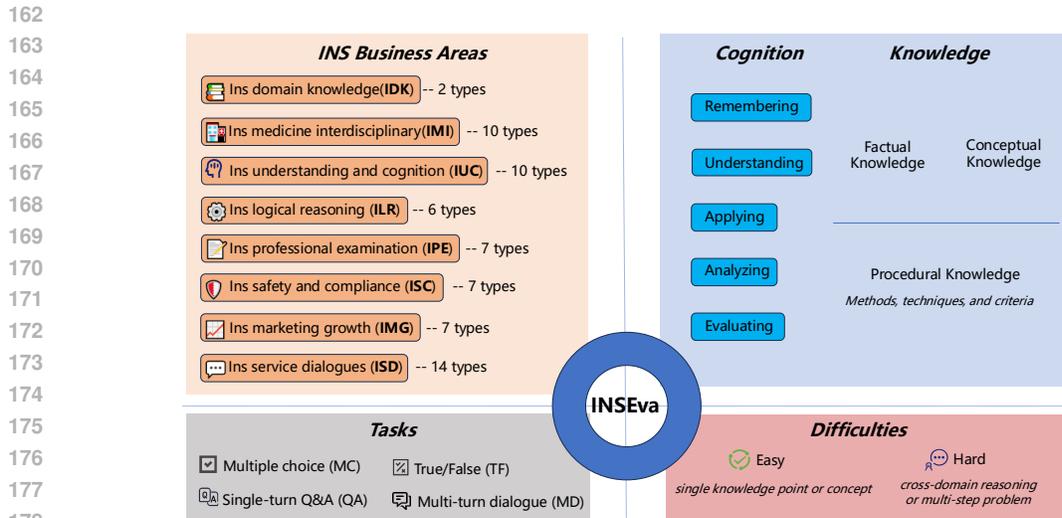


Figure 2: The taxonomy of INSEva benchmark.

3.2 TAXONOMY

To comprehensively assess models’ capabilities in the insurance domain, we propose a multi-dimensional taxonomy for our evaluation benchmark. This taxonomy categorizes test items along four primary dimensions, including the business areas dimension, the task formats dimension, the difficulties dimension, and the cognition and knowledge dimension. The multi-dimensional taxonomy enables a fine-grained analysis of model performance across different aspects of insurance knowledge and applications.

3.2.1 DIMENSION 1: BUSINESS AREAS

We define a three-pillar, eight-dimension task taxonomy by integrating three sources of requirements: (i) the complementary strengths of general LLMs versus domain LLMs (language understanding/reasoning/safety alignment v.s. deep domain knowledge); (ii) the competency framework reflected in undergraduate and graduate insurance curricula (e.g., risk and actuarial science, insurance law and regulation, health/medical knowledge for health insurance, marketing and customer operations); and (iii) capability demands observed in end-to-end insurance applications (underwriting, claims, compliance, sales, and service). This leads to three capability pillars—Knowledge Literacy, Foundational General Skills, and Business-Integrated Skills—instantiated by eight business-area dimensions that jointly cover domain knowledge, cognition and reasoning, compliance, and workflow-centric dialogue and generation.

We categorize examples as follows: (1) **Insurance Domain Knowledge (IDK)** assesses professional knowledge literacy via two tasks—Insurance Knowledge Interpretation and Insurance Science; (2) **Insurance-Medical Interdisciplinary (IMI)** comprises ten cross-domain tasks spanning medical entity extraction and standardization, diagnosis/Q&A (including pet scenarios), and risk/prescription prediction; (3) **Insurance Understanding & Cognition (IUC)** includes ten information extraction and analysis tasks such as intent understanding, slot filling for products/insured objects/diseases, attribute extraction, clause interpretation, and liability/product selection analysis; (4) **Insurance Logical Reasoning (ILR)** emphasizes six quantitative and rule-based reasoning tasks, including actuarial science, financial mathematics, numerical computation, and exemption reasoning; (5) **Insurance Professional Examination (IPE)** simulate seven certification settings across insurance, physician/pharmacist/veterinarian practice, actuarial qualification, and salesperson licensing; (6) **Insurance Safety & Compliance (ISC)** defines seven tasks on information security, baseline controls, document compliance, value alignment, issue identification, fact-checking, and compliance verification; (7) **Insurance Marketing Growth (IMG)** evaluates seven tasks on customer segmentation, service summarization, marketing copy generation, recommendation script-

ing, investor education, population classification, and strategy formulation; (8) **Insurance Service Dialogue (ISD)** constructs fourteen multi-turn tasks covering the service lifecycle, including product/regulatory interpretation, underwriting and policy servicing, claims assessment and settlement, post-policy operations, planning/configuration, condition-based selection and comparison, and premium/benefit calculation. More details are provided in the Appendix.

3.2.2 DIMENSION 2: TASK FORMATS

According to the task format of the example, we classify examples into four types: (1) **Multiple choice (MC)**: Tasks requiring selection of the correct option(s) from a set of alternatives; (2) **True/False (TF)**: Binary judgment tasks requiring verification of statements' accuracy; (3) **Single-turn Q&A (QA)**: Direct question-answer pairs requiring concise responses; (4) **Multi-turn dialogue (MD)**: Extended conversations simulating real-world insurance consultations requiring contextual understanding and knowledge application across multiple exchanges.

3.2.3 DIMENSION 3: DIFFICULTIES

To assess the scalability of models across varying levels of expertise required for insurance operations, we categorize examples according to their difficulty: (1) **Easy**: Items can be answered directly without requiring deep understanding, or items only requiring the retrieval and application of a single knowledge point or concept; (2) **Hard**: Complex scenarios requiring integration and application of multiple knowledge elements, often involving cross-domain reasoning or multi-step problem-solving.

3.2.4 DIMENSION 4: COGNITION & KNOWLEDGE

To enable a granular assessment of models' cognitive processes and knowledge within the insurance domain, informed by established education assessment theory (Krathwohl, 2002), we categorize examples according to their alignment with specific cognitive dimensions and knowledge dimensions¹.

For the cognition dimension: (1) **Remembering (Rem.)**: Recalling or recognizing basic insurance knowledge (e.g., facts, terminology) without deeper understanding; (2) **Understanding (Und.)**: Comprehending insurance concepts through explanation, transformation, or inference to establish basic connections; (3) **Applying (App.)**: Using acquired insurance knowledge or methods to solve problems in new contexts; (4) **Analyzing (Ana.)**: Breaking down insurance information to clarify relationships, structures, and theories; (5) **Evaluating (Eva.)**: Critically assessing the logic, structure, or effectiveness of insurance information and ideas.

For the knowledge dimension: (1) **Factual Knowledge (FK)**: Basic elements of insurance (e.g., terminology, specific details); (2) **Conceptual Knowledge (CK)**: Interrelationships among basic elements within insurance theory (e.g., classifications, principles); (3) **Procedural Knowledge (PK)**: Methods, techniques, and criteria for determining when to use appropriate insurance procedures.

3.3 DATA CONSTRUCTION

This section details our systematic approach to developing a comprehensive benchmark in the insurance domain, encompassing data collection, data augmentation, and data quality control methods.

3.3.1 DATA COLLECTION

Our benchmark integrates data from diverse authoritative sources to ensure comprehensiveness and domain relevance. (1) We systematically extract historical examination questions, corresponding answers, and explanatory notes from professional Chinese certification platforms, subsequently restructuring these data to conform to the format; (2) We collect specialized terminology and regulatory standards from authoritative industry repositories, from which we extract core knowledge points and reformulate them into question-answer formats using LLM; (3) We curate internal business resources containing domain-specific knowledge and scenarios, transforming them into standardized

¹It is worth noting that in the education assessment theory, there is also Creating under the cognition dimension and Metacognitive Knowledge under the knowledge dimension. Due to the seriousness and professionalism of the insurance domain, these two dimensions are not included in our taxonomy.

Table 1: Data statistics of different business areas and cognition and knowledge in INSEva. The absolute counts and relative percentage distribution are presented. The abbreviations are defined in Section 3.2.

Business	Tasks	Len.	Num.	%
IDK	MC	226	1336	3.5
IMI	MC/TF	372	9447	24.4
IUC	MC	402	7916	20.5
ILR	MC/TF	603	3923	10.1
IPE	MC	195	7932	20.5
ISC	TF	169	3417	8.1
IMG	MC	1141	1651	4.3
ISD	QA/MD	3230	3082	8.0

(a) Business Areas

Type	Tags	Num.	%
Cognition	Rem.	541	1.4
	Und.	8036	20.8
	App.	19339	50.0
	Ana.	7038	18.2
Knowledge	Eva.	3750	9.7
	FK	914	2.4
	CK	21047	54.4
	PK	16743	43.3

(b) Cognition and Knowledge

question-answer pairs using LLM. Finally, we collect about 40k examples from different sources. This multi-source approach ensures our benchmark encompasses both theoretical knowledge and practical applications within the insurance domain.

3.3.2 DATA AUGMENTATION

To enhance the robustness and validity of our benchmark, we implement several augmentation strategies for questions and answers on the collected data. (1) **Question Augmentation**: To promote linguistic diversity, we employ paraphrasing techniques to generate semantically equivalent questions with distinct syntactic expressions, thereby reducing potential model bias toward specific phrasings. Additionally, we make colloquial rewrites for standard questions, converting formal insurance terminology into colloquial expressions more representative of real-world customer interactions. This transformation enhances the benchmark’s utility for customer-facing applications; (2) **Answer Augmentation**: We implement option randomization to mitigate selection bias (Zheng et al., 2024) for multiple-choice questions, ensuring uniform distribution of correct answers across options throughout the dataset. Finally, we get nearly 80k augmented examples in the benchmark.

3.3.3 DATA QUALITY CONTROL

To ensure the reliability and validity of our benchmark, we implement a multi-faceted quality control framework comprising rule-based, expert-based, and LLM-based validation modules. (1) **Rule-Based Module**: This module evaluates the quality of the samples through a series of rules. Specifically, we conduct: (1) character set validation to identify and rectify illegal characters, special symbols, or encoding anomalies, thereby ensuring Unicode compliance and eliminating corrupted elements; (2) structural validity verification for multiple-choice questions, confirming option completeness, absence of duplicates, and adherence to specified formatting requirements; and (3) content-specific quality checks, wherein samples are assessed against pre-defined criteria based on their inherent properties (e.g., questions about insurance products must explicitly reference the product name); (2) **Expert-Based Module**: This module ensures the faithfulness and consistency of ground truth in Q&A and dialogue samples through rigorous domain expert examination. A team of 10 insurance industry experts, averaging five years of experience and proficient in both Chinese and English, independently reviews the data to verify that ground truth responses are fully supported by and logically consistent with the provided input context. Any inconsistent or unsupported examples are removed. Inter-annotator agreement, measured by Cohen’s Kappa, reaches 0.87, reflecting high consistency. Detailed annotation guidelines and tools are employed to ensure uniformity and reproducibility, as shown in Appendix E; (3) **LLM-Based Module**: To further enhance validation reliability and mitigate potential biases, we employ a multi-LLM voting approach. In this module, multiple distinct LLMs independently evaluate each sample. Final quality determinations are then made through aggregation of these independent assessments, leveraging the diverse perspectives of multiple models to achieve a more robust and objective validation process.

3.4 DATA STATISTICS

Finally, we attain a total of 38,704 examples, while the average length of the prompts of these examples is 905 tokens. Table 1a and Table 1b show more details of different business areas and different cognition and knowledge. A substantial proportion of questions in the benchmark belong to knowledge applying and procedural knowledge types, which provide a more realistic assessment of the model’s practical capabilities in insurance.

3.5 EVALUATION METHODS

Due to the high risks of the insurance domain, we pay more attention to the correctness of the model’s response in the insurance domain. Our evaluation framework employs different assessment approaches tailored to various question types in the benchmark.

For deterministic questions, such as multiple-choice and true/false questions, whether the model’s response is the same as the ground-truth represents the correctness of the model’s response. Specifically, we calculate *Accuracy* as the primary metric, representing the proportion of questions for which the model selects the correct option from the available choices. We further stratify performance analysis across different dimensions of our taxonomy to identify specific strengths and weaknesses in model capabilities.

For open-ended questions, including single-turn Q&A and multi-turn dialogues, we not only require that the model’s response cover the ground truth as much as possible, but also require that no hallucinations appear in the response. Therefore, we propose two metrics, *faithfulness* and *completeness*, and implement an LLM-based evaluation pipeline for each metric that approximates human expert assessment. *Faithfulness* aims to measure the consistency of a model’s response with the provided context, reflecting whether the answer contains hallucinations or unsupported claims. To calculate the *faithfulness*, we decompose the model’s response into individual statements and determine whether the provided retrieval context substantiates each statement. Specifically, we use an LLM to decompose the model’s response into a maximum of 20 statements. For each statement, the LLM determines whether it is supported by the provided retrieval context, assigning a binary score of 1 if supported and 0 otherwise. The overall *faithfulness* score is then calculated as the proportion of supported statements. *Completeness* aims to measure the extent to which a model’s response covers the ground truth, reflecting whether the answer covers the valid information required. The core method involves decomposing the ground truth into statements and assessing whether each statement is present in the model’s response. Specifically, we employ an LLM to decompose the ground truth into a maximum of 20 statements. The LLM then determines whether each statement is present within the model’s response, assigning a binary score of 1 if present and 0 otherwise. The overall *Completeness* score is computed as the proportion of ground truth statements that are successfully recalled in the model’s response.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

We select 9 state-of-the-art LLMs for evaluation, including closed-source, open-source, and domain-specific models, which are as follows: (1) **Closed-source:** GPT-4o (OpenAI, 2024); Doubao-1.5-pro-256k²; Gemini-2.5-pro³; (2) **Open-source:** Qwen2.5-72B-Instruct (Qwen et al., 2025); Qwen-QwQ-32B (Team, 2025); Deepseek-R1 (Guo et al., 2025); Qwen3-235B-A22B (Yang et al., 2025); (3) **Domain-specific:** Fin-R1 (Liu et al., 2025); DianJin-R1 (Zhu et al., 2025).

4.2 MAIN RESULTS

Based on pre-defined taxonomy and evaluation methods, we conduct a comprehensive assessment, and the results are presented in Table 2. A comprehensive evaluation of LLMs across different dimensions specific to the insurance domain reveals significant variations in overall performance.

²https://seed.bytedance.com/en/special/doubao_1_5_pro

³<https://deepmind.google/models/gemini/pro/>

Table 2: Main results for LLMs on INSEva across different business areas (abbreviations defined in Section 3.2).

Models	IDK	IMI	IUC	ILR	IPE	ISC	IMG	ISD		Avg.
								Faithful.	Complete.	
GPT-4o	84.24	79.25	84.90	56.68	74.83	73.84	94.85	82.26	80.97	79.09
Doubao-1.5	90.17	80.88	85.44	58.75	87.11	78.35	94.93	86.86	82.17	82.74
Gemini-2.5	88.14	80.68	84.54	78.27	87.75	79.13	95.75	75.54	85.24	83.89
Qwen2.5	86.95	77.79	82.90	65.72	80.23	80.45	95.17	78.99	85.36	81.51
Qwen-QwQ	89.05	80.07	84.62	70.90	81.11	76.86	94.98	73.94	89.47	82.33
Deepseek-R1	89.48	79.75	85.00	73.66	86.15	70.65	95.71	74.51	87.76	82.52
Qwen3	86.48	79.51	84.62	71.74	85.01	73.34	95.29	72.87	88.40	81.92
Fin-R1	80.34	72.48	76.40	61.90	67.20	74.56	87.41	70.70	81.04	74.67
DianJin-R1	87.23	72.59	82.80	69.36	82.65	70.19	85.58	71.21	87.55	78.80

Gemini-2.5-pro demonstrates the highest aggregate performance, achieving a mean score of 83.89, while the domain-specific Fin-R1 model exhibits a comparatively lower mean score of 74.67. Notably, with the exception of the domain-specific model, the general LLMs (such as Qwen2.5, Deepseek-R1, etc.) maintain average scores above 80, indicating a relatively stable foundational capability in the insurance domain. However, a gap remains between current LLM performance and the level of expertise required to fully address complex real-world insurance problems.

Regarding granular dimensions, all evaluated models achieve excellent results in the Insurance Marketing Growth (IMG), with an average score of more than 94 points. The reason is that although IMG is composed of insurance data, the tasks in it are similar to general tasks in NLP (such as text summarization, etc.), and require the general language understanding and generation capabilities of LLMs. However, they generally perform poorly in Insurance Logical Reasoning (ILR), with scores generally lower than other dimensions, where the highest score is only 78.27 of Gemini-2.5-pro. This uneven distribution of capabilities reflects that the current LLMs still have room for improvement in complex reasoning tasks. Additionally, the generally low scores in the Insurance Safety Compliance (ISC) also expose the limitations of the model in dealing with strong regulatory scenarios. The mediocre performance of the models in these domain-specific tasks further emphasizes the need for insurance benchmarks to fill the gap in capability evaluation in the insurance domain.

Different LLMs show obvious differentiated capabilities. Gemini-2.5-pro achieves competitive performance in all dimensions, showing strong general performance. Qwen series LLMs have relative advantages in reasoning and dialogue generation. Fin-R1 and DianJin-R1 show typical domain-specific characteristics. They have excellent performance in the financial benchmark, but their performance in insurance is even worse than that of general LLMs. It shows that even though insurance is a subdomain of finance and there are many existing financial benchmarks, it is still very important to construct an insurance benchmark.

Through an in-depth analysis of the faithfulness and completeness under the insurance service dialogue (ISD), we can find that different types of LLMs show significant trade-off characteristics in these two key indicators. Specifically, reasoning models represented by Deepseek-R1 perform well in completeness, but their faithfulness index is relatively low. It shows that this type of model tends to generate more comprehensive answers through reasoning and knowledge integration, which can better cover the key information points in the ground truth, but at the same time it is also more likely to generate additional content that is not explicitly mentioned in the dataset. In contrast, the non-reasoning LLMs represented by GPT-4o show relatively high faithfulness, but the completeness is relatively low, reflecting that this type of model is more inclined to generate strict answers based on known information. This trade-off relationship between faithfulness and completeness has special practical significance in high-risk financial fields such as insurance. Insurance products usually involve complex clause details and strict regulatory requirements. False or inaccurate information generated by the model may cause users to make wrong insurance decisions, leading to serious economic losses or even legal disputes. Therefore, in practical applications, we believe that faithfulness should be given priority to ensure the reliability of model output.

Main results for LLMs across different cognition and knowledge are shown in Table 3. The experimental results reveal several noteworthy patterns across cognitive and knowledge dimensions. In terms of cognition dimensions, all models demonstrate stronger performance in lower-order cogni-

Table 3: Main results for LLMs across different cognition and knowledge.

Models	Cognition					Knowledge			Avg.
	Rem.	Und.	App.	Ana.	Eva.	FK	CK	PK	
GPT-4o	80.30	87.74	80.67	74.95	71.33	83.57	82.14	77.13	79.72
Doubao-1.5	89.01	88.93	86.14	77.12	72.00	88.67	87.29	78.35	83.43
Deepseek-R1	91.39	89.28	87.30	81.33	69.24	91.57	84.44	83.58	84.76
Qwen3	88.23	87.81	85.89	81.71	75.17	88.69	84.82	83.59	84.48
DianJin-R1	85.02	85.51	81.12	73.80	69.36	82.83	81.43	76.47	79.44

Table 4: Experimental results comparing insurance-specific model Finix-S1 with the base model.

Models	IDK	IMI	IUC	ILR	IPE	ISC	IMG	ISD		Avg.
								Faithful.	Complete.	
Qwen-QwQ	89.05	80.07	84.62	70.90	81.11	76.86	94.98	73.94	89.47	82.33
Finix-S1	92.66	91.13	89.69	78.69	87.15	83.49	90.07	88.81	81.45	87.02
Δ	+3.61	+11.06	+5.07	+7.79	+6.04	+6.63	-4.91	+14.87	-8.02	+4.69

tive tasks such as Remembering (Rem.) and Understanding (Und.), with scores consistently above 85%, while showing relative weakness in higher-order cognitive tasks, particularly in Evaluation (Eva.) where scores drop below 75%. Notably, Deepseek-R1 exhibits exceptional performance in basic cognitive tasks (91.39% in Remembering), but experiences a significant performance degradation in evaluation tasks (69.24%), suggesting a common challenge in high-level cognition across all models. In terms of knowledge dimensions, models generally perform better in Factual Knowledge (FK) compared to Procedural Knowledge (PK), with Deepseek-R1 achieving the highest FK score of 91.57%. This pattern indicates that current models have a stronger grasp of basic facts and terminology than procedural operations in the insurance domain. The overall performance comparison shows Deepseek-R1 and Qwen3 leading with average scores of 84.76% and 84.48% respectively, significantly outperforming GPT-4o (79.72%) and DianJin-R1 (79.44%). This performance gap suggests that recent architectural improvements and training strategies have effectively enhanced models’ capabilities across both cognition and knowledge dimensions.

4.3 INSURANCE-SPECIFIC MODEL TRAINING AND EVALUATION

To validate the effectiveness of the data construction pipeline, we also train a domain-specific model named **Finix-S1** based on Qwen-QwQ-32B model using reinforcement learning. As shown in Table 4, evaluation results demonstrate that Finix-S1 achieves a significant performance improvement. It surpasses the base model Qwen-QwQ by 4.69% on average, and current state-of-the-art general LLMs (e.g., Gemini) by over 3% on average. This advantage is particularly pronounced on the Insurance-Medical Interdisciplinary (IMI), which necessitates the integration of knowledge across multiple domains, with Finix-S1 exceeding the best model by more than 10%. This suggests that prevailing general LLMs exhibit limitations in cross-domain knowledge integration, highlighting the need for training paradigms that foster this capability. However, a key trade-off is observed: on tasks akin to general NLP (e.g., IMG), Finix-S1’s performance lags behind that of general LLMs. This phenomenon, consistent across other domain-specific models, points to the challenge of catastrophic forgetting, where gains in specialized performance may come at the expense of general capabilities. More training and evaluation details are provided in the Appendix C.

5 CONCLUSION

We introduce INSEva, a comprehensive benchmark for the insurance domain featuring 38,704 high-quality examples. Our evaluations reveal that existing LLMs, including those specialized in finance, possess only basic competency and struggle with complex reasoning, often trading faithfulness for completeness. This highlights the unique challenges of the insurance sector. Ultimately, INSEva serves as a vital tool for practitioners in model selection and guides researchers in developing more robust architectures for high-risk domains.

6 ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. In this study, no human subjects or animal experimentation was involved. All datasets used were sourced in compliance with relevant usage guidelines, ensuring no violation of privacy. We have taken care to avoid any biases or discriminatory outcomes in our research process. No personally identifiable information was used, and no experiments were conducted that could raise privacy or security concerns. We are committed to maintaining transparency and integrity throughout the research process.

7 REPRODUCIBILITY STATEMENT

We are committed to fostering reproducible research. To this end, all components of our evaluation framework, including detailed descriptions of the methodology, evaluation scripts, and model interaction protocols, will be open-sourced upon publication. This will enable other researchers to apply our evaluation approach to their own models.

Regarding the dataset, the full benchmark contains sensitive commercial information and is proprietary, which prevents its public release. To balance the need for transparency with our confidentiality obligations, we will make a 10% stratified sample of the benchmark publicly accessible. This sample is carefully curated to reflect the distribution and complexity of the complete benchmark, providing a valuable resource for future research and analysis.

REFERENCES

- Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, et al. Finqa: A dataset of numerical reasoning over financial data. *arXiv preprint arXiv:2109.00122*, 2021.
- Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering. *arXiv preprint arXiv:2210.03849*, 2022.
- Jing Ding, Kai Feng, Binbin Lin, Jiarui Cai, Qiushi Wang, Yu Xie, Xiaojin Zhang, Zhongyu Wei, and Wei Chen. Insqabench: Benchmarking chinese insurance domain question answering with large language models, 2025. URL <https://arxiv.org/abs/2501.10943>.
- Minwei Feng, Bing Xiang, Michael R. Glass, Lidan Wang, and Bowen Zhou. Applying deep learning to answer selection: A study and an open task, 2015. URL <https://arxiv.org/abs/1508.01585>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>.
- Ravi Jaiswal. Impact of ai in the general insurance underwriting factors. *Central European Management Journal*, pp. 697–705, 2023.
- David R Krathwohl. A revision of bloom’s taxonomy: An overview. *Theory into practice*, 41(4): 212–218, 2002.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekogun, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori

- 540 Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yi-
541 fan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models, 2023. URL
542 <https://arxiv.org/abs/2211.09110>.
- 543
544 Chenwei Lin, Hanjia Lyu, Xian Xu, and Jiebo Luo. Ins-mmbench: A comprehensive benchmark for
545 evaluating lvlms’ performance in insurance, 2024. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.09105)
546 09105.
- 547 Zhaowei Liu, Xin Guo, Fangqi Lou, Lingfeng Zeng, Jinyi Niu, Zixuan Wang, Jiajie Xu, Weige Cai,
548 Ziwei Yang, Xueqian Zhao, Chao Li, Sheng Xu, Dezhi Chen, Yun Chen, Zuo Bai, and Liwen
549 Zhang. Fin-rl: A large language model for financial reasoning through reinforcement learning,
550 2025. URL <https://arxiv.org/abs/2503.16252>.
- 551
552 OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- 553 Xueqing Peng, Triantafillos Papadopoulos, Efstathia Soufleri, Polydoros Giannouris, Ruoyu Xi-
554 ang, Yan Wang, Lingfei Qian, Jimin Huang, Qianqian Xie, and Sophia Ananiadou. Plutus:
555 Benchmarking large language models in low-resource greek finance, 2025. URL <https://arxiv.org/abs/2502.18772>.
- 556
557 Robert Pfeifer and Franziska Langen. Insurance and sustainable development goals. 2021. <https://arxiv.org/pdf/2102.02612>.
- 558
559
560 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan
561 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,
562 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin
563 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,
564 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,
565 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.
566 URL <https://arxiv.org/abs/2412.15115>.
- 567 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
568 Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof qa bench-
569 mark, 2023. URL <https://arxiv.org/abs/2311.12022>.
- 570
571 Raj Sanjay Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj
572 Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. When flue meets flang: Benchmarks and
573 large pre-trained language model for financial domain. *arXiv preprint arXiv:2211.00083*, 2022.
- 574 Rao A et al. Srivastava A, Rastogi A. Beyond the imitation game: Quantifying and extrapolating the
575 capabilities of language models, 2023. URL <https://arxiv.org/abs/2206.04615>.
- 576
577 Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL
578 <https://qwenlm.github.io/blog/qwq-32b/>.
- 579 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
580 Glue: A multi-task benchmark and analysis platform for natural language understanding, 2019.
581 URL <https://arxiv.org/abs/1804.07461>.
- 582
583 Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao,
584 Dong Li, Yongfu Dai, Duanyu Feng, et al. Finben: A holistic financial benchmark for large
585 language models. *Advances in Neural Information Processing Systems*, 37:95716–95743, 2024.
- 586 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
587 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
588 *arXiv:2505.09388*, 2025.
- 589
590 Liwen Zhang, Weige Cai, Zhaowei Liu, Zhi Yang, Wei Dai, Yujie Liao, Qianru Qin, Yifei Li, Xingyu
591 Liu, Zhiqiang Liu, et al. Fineval: A chinese financial domain knowledge evaluation benchmark
592 for large language models. *arXiv preprint arXiv:2308.09975*, 2023.
- 593 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are
not robust multiple choice selectors, 2024. URL <https://arxiv.org/abs/2309.03882>.

Table 5: Cross-Lingual evaluation results. The performance of several LLMs on a small constructed English benchmark is similar to that on Chinese.

Models	IDK	IMI	IUC	ILR	IPE	ISC	IMG	ISD		Avg.
								Faithful.	Complete.	
Doubao-1.5	86.67	79.38	82.61	59.16	74.12	71.15	90.00	84.45	80.62	78.68
Qwen2.5	84.59	75.85	81.62	64.73	68.50	92.82	94.14	85.51	71.15	79.88
Qwen-QwQ	86.67	77.57	83.24	70.82	75.29	81.66	93.09	70.48	83.26	80.23
Qwen3	85.84	76.71	81.88	75.70	72.07	73.99	93.93	75.73	80.13	79.55
Fin-R1	69.58	67.56	78.15	71.03	57.05	89.96	83.69	62.28	72.23	72.39

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. Tat-qa: A question answering benchmark on a hybrid of tabular and textual content in finance, 2021. URL <https://arxiv.org/abs/2105.07624>.

Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. Benchmarking large language models on cflue - a chinese financial language understanding evaluation dataset. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics(ACL-2024)*, 2024.

Jie Zhu, Qian Chen, Huaixia Dou, Junhui Li, Lifan Guo, Feng Chen, and Chi Zhang. Dianjin-r1: Evaluating and enhancing financial reasoning in large language models. *arXiv preprint arXiv:2504.15716*, 2025.

A USE OF LARGE LANGUAGE MODELS

During the preparation of this work, we used LLMs solely for the purpose of improving language and clarity. Specifically, LLMs were used for proofreading, grammar correction, and minor phrasing improvements.

We reviewed and edited all output generated by the LLMs, and take full responsibility for all content and ideas presented in this work.

B EXPERIMENTS

B.1 CROSS-LINGUAL EVALUATION

We conduct several cross-lingual experiments, as shown in Table 5. Specifically, we sample approximately 30% of the Chinese questions (around 10k) from our benchmark and translate them into English, and then test them across several models. Although the overall performance on English is slightly lower, the key conclusions remain consistent: performance is relatively low on ILR, a trade-off exists between faithfulness and completeness, and domain-specific models underperformed compared to general models. The results demonstrate that the benchmark is designed to assess the domain knowledge and reasoning abilities of LLMs in insurance, which remain consistent across languages without strong linguistic dependency.

B.2 CORRELATION ANALYSIS

To validate the effectiveness of evaluation methods for open-ended questions, we select 98 instances from different categories in the INSEva benchmark and choose 4 different models, totally 392 responses for correlation analysis. We choose the BLEU and ROUGE family of methods as our baselines, which are widely used metrics for text quality assessment. Human evaluators are asked to rate the models’ responses on faithfulness and completeness. The human-rated scores are considered as the optimal evaluation method, and we conduct a correlation analysis between various common automatic evaluation methods and the proposed evaluation method.

The correlation results are presented in Table 6. It can be observed that our evaluation method exhibits the best correlation, which means our evaluation method can achieve high consistency with human evaluation.

Table 6: Sample-level correlation of different metrics.

Metrics	Faithful.	Complete.	Avg.
Bleu-1	31.47	-18.85	6.31
Bleu-2	37.05	-15.32	10.87
Bleu-4	44.01	-6.07	18.97
Rouge-1	38.71	-17.45	10.63
Rouge-2	49.53	-1.06	24.24
Rouge-L	43.86	-7.63	18.12
Ours	73.60	55.78	64.69

Table 7: Experimental results for Finix-S1 with varying amounts of insurance-specific training data, while holding general data fixed.

Models	IDK	IMI	IUC	ILR	IPE	ISC	IMG	ISD		Avg.
								Faithful.	Complete.	
Qwen-QwQ	89.05	80.07	84.62	70.90	81.11	76.86	94.98	73.94	89.47	82.33
0%	88.74	79.31	82.57	70.81	78.04	75.09	94.75	75.63	86.02	81.22
50%	90.53	87.47	87.74	76.91	85.48	80.97	90.93	82.13	84.33	85.16
75%	91.87	90.47	88.83	78.03	86.74	82.53	90.03	86.32	82.17	86.33
100%	92.66	91.13	89.69	78.69	87.15	83.49	90.07	88.81	81.45	87.02

C DETAILS OF FINIX-S1

Finix-S1’s training consists of two stages, conducted on NVIDIA A100 GPUs. The model is first trained via Supervised Fine-Tuning (SFT) using 32 GPUs across 4 nodes. Then, we employ reinforcement learning to implement Group Relative Policy Optimization (GRPO), leveraging insurance Q&A data and incorporating a dual reward mechanism to improve both the accuracy of response formatting and content. In the RL stage, we perform 3 rollouts per sample, with a train batch size of 72.

The composition of our training data is as follows: (1) 108K general data samples collected from the GPQA dataset (Rein et al., 2023); (2) 347.3K samples constructed following our benchmark construction pipeline for SFT, along with an additional 36.5K samples utilized for reinforcement learning; (3) All training data are strictly separated from the evaluation data, ensuring that there is no data leakage between the training and evaluation sets.

We also evaluate the model by fixing the amount of general training data and varying the quantity of insurance-specific training data, as shown in Table 7. The results demonstrate that increasing the amount of insurance-specific training data consistently improves performance across multiple business areas. This indicates that current LLMs indeed lack sufficient insurance domain knowledge and the ability to apply knowledge. Furthermore, these results validate the effectiveness of our insurance data construction methodology. However, a slight performance decrease was observed specifically on the IMG. As discussed in the experimental section, this minor degradation aligns with the phenomenon of catastrophic forgetting, where increased specialization in the insurance domain leads to a decline in performance on certain more general tasks.

Here is an example. The user asks whether “hospitalization due to trauma within the past two years” can be covered by the Hao Yi Bao Long-term Medical Insurance (20-year edition). According to the standards, for first-time or non-continuous applicants, any illness or injury occurring before the inception of the insurance contract (categorized as a pre-existing condition, including trauma) is excluded from coverage. Finix-S1 with 0% insurance training data generates an erroneous response: “The system asserts that trauma is classified as accidental injury and, if the insurance is renewed continuously and no waiting period applies, coverage might be possible. It concludes with possible compensation after discussing different scenarios.” While Finix-S1 with 100% insurance training data generates a correct response: “Hospitalization resulting from trauma within the past two years constitutes a pre-existing condition. According to the policy, this is categorically excluded from coverage (denied compensation).” The inadequately trained model demonstrates a **misunderstanding of domain-specific terminology**, specifically the distinction between “waiting period” and

702 “pre-existing condition” in insurance contracts. By conflating the rules for waiting periods—which
703 primarily apply to post-contract incidents—with the distinct and absolute nature of pre-existing con-
704 dition exclusions, the model erroneously infers that continuous renewal could eliminate the exclu-
705 sion for prior trauma. This indicates a lack of nuanced comprehension of insurance lexicon and the
706 relationships between contract clauses, leading to an incorrect compensation assessment for prior
707 hospitalizations.

708 709 D DETAILS OF TAXONOMY 710

711 Our Insurance Business Areas dimensions consists of eight main areas:
712

- 713 1. Insurance Domain Knowledge (IDK):
 - 714 • Insurance Knowledge Interpretation
 - 715 • Insurance Science
- 716 2. Insurance-Medicine Interdisciplinary (IMI):
 - 717 • Medical Procedure Name Extraction
 - 718 • Medical Condition Name Extraction
 - 719 • Disease Relevance Discrimination
 - 720 • Insurance Hospital Name Standardization
 - 721 • Insurance Disease & Procedure Standardization
 - 722 • Pet Disease Diagnosis
 - 723 • Pet Disease Examination Q&A
 - 724 • Pet Medication Knowledge Q&A
 - 725 • Medical Prescription Disease Prediction
 - 726 • Disease Risk Prediction in Medical Reports
- 727 3. Insurance Understanding and Cognition (IUC):
 - 728 • Insurance Intent Understanding
 - 729 • Insurance Product Slot Recognition
 - 730 • Insured Object Slot Recognition
 - 731 • Insurance Disease Slot Recognition
 - 732 • Insurance Attribute Extraction
 - 733 • Insurance Clause Interpretation
 - 734 • Insurance Product Selection Analysis
 - 735 • Insurance Liability Analysis
 - 736 • Insurance Review Tag Recognition
 - 737 • Insurance Product Review Classification
- 738 4. Insurance Logical Reasoning (ILR):
 - 739 • Insurance Actuarial Science
 - 740 • Financial Mathematics
 - 741 • Financial Numerical Computation
 - 742 • Insurance Prior Exemption Reasoning
 - 743 • Insurance General Exemption Reasoning
 - 744 • Pure Outpatient Insurance Reasoning
- 745 5. Insurance Professional Examination (IPE):
 - 746 • Insurance Professional Qualification Exams
 - 747 • CICE Insurance Certification
 - 748 • Practicing Physician Qualification Exam
 - 749 • Practicing Pharmacist Qualification Exam
 - 750 • Practicing Veterinarian Qualification Exam
 - 751 • Chinese Actuary Examination

- 756 • Insurance Salesperson Certification Exam
- 757
- 758 6. Insurance Safety and Compliance (ISC):
- 759 • Information Security
- 760 • Security Baseline
- 761 • Insurance Document Compliance
- 762 • Insurance Value System
- 763 • Insurance Issue Identification
- 764 • Insurance Fact-Checking
- 765 • Insurance Compliance Verification
- 766
- 767
- 768 7. Insurance Marketing Growth (IMG):
- 769 • Insurance Target Population Positioning
- 770 • Insurance Service Summary
- 771 • Insurance Marketing Copy Generation
- 772 • Insurance Product Recommendation Scripting
- 773 • Insurance Investor Education Scripting
- 774 • Insurance Population Identification & Classification
- 775 – User Purchase Intention Classification
- 776 – User Cognitive Level Classification
- 777 – User Attitude Level
- 778 – Insurance Type Preference
- 779
- 780 • Insurance Service Strategy Formulation
- 781 – Insurance Service Necessity Determination
- 782 – Insurance Service Scenario Selection
- 783 – Insurance Service Timing Decision
- 784 – Insurance Service Scenario Expression
- 785
- 786
- 787 8. Insurance Service Dialogues (ISD):
- 788 • Pre-Investment Consultation
- 789 – Product Interpretation
- 790 – Regulatory Information Interpretation
- 791 – Underwriting-related
- 792 – Policy-related
- 793 – Platform-related
- 794
- 795 • Post-Investment Issue Resolution
- 796 – Claims Assessment-related
- 797 – Claim Settlement-related
- 798 – Post-policy Operations
- 799
- 800 • Professional Allocation Advice
- 801 – Planning and Configuration
- 802 – Condition-based Product Selection
- 803 – Product Recommendation
- 804 – Insurance Type Comparison
- 805 – Product Comparison
- 806 – Premium and Benefit Calculation
- 807

808 This comprehensive taxonomy covers the full spectrum of insurance business areas, from funda-
809 mental domain knowledge to practical service applications, ensuring a thorough evaluation of AI
models in the insurance domain.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

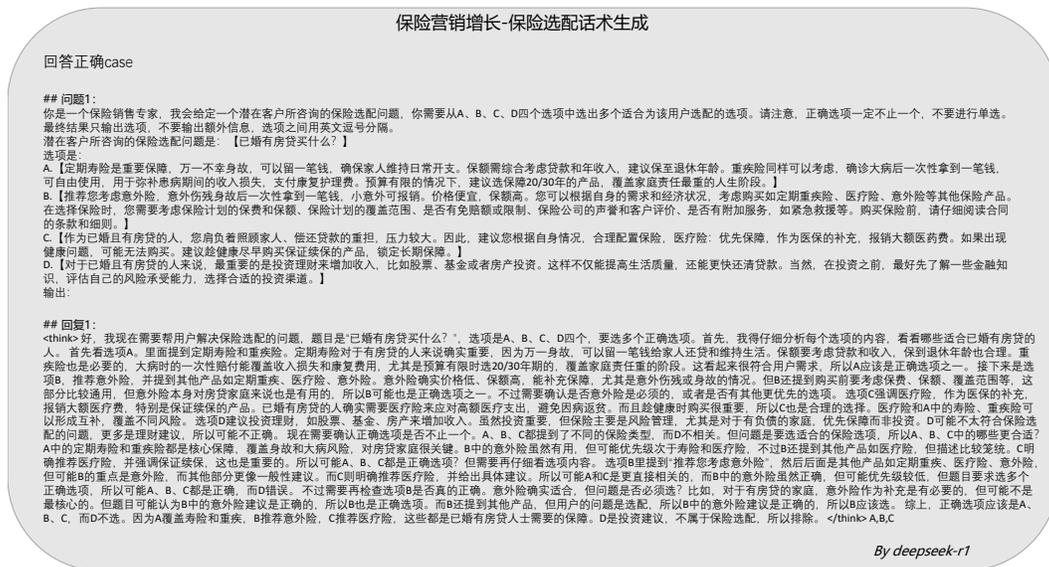


Figure 3: A case of correct answer in Chinese.

E DETAILS OF HUMAN ANNOTATIONS

The human annotations in the domain expert module and correlation analysis consist of over 10 insurance industry experts with an average of 5 years of experience, based in China. They are proficient in both Chinese and English, belong to different organizations respectively, and can utilize their industry experience along with professional reference materials to analyze and evaluate answers to the questions.

Human experts have established annotation criteria that primarily evaluate answer quality from three dimensions - factuality, accuracy, and continuity - through granular assessment of knowledge points. The evaluation scale ranges from 5 (excellent) to 0 (poor):

- **Factuality:** Whether the response content adheres to factual truth and can be substantiated by evidence from retrieved knowledge sources.
- **Accuracy:** Whether the response demonstrates professional rigor without containing errors in knowledge, expertise, or factual information.
- **Continuity:** Whether the response maintains contextual coherence, enables smooth communication flow, and exhibits clear language expression through effective discourse transitions.

F CASE STUDY

As illustrated in Figure 3 and Figure 4, DeepSeek-R1 leverages step-by-step deliberation to correctly select the life-critical triad of term-life, critical-illness and medical coverage for a mortgaged family, evidencing its strength in multi-hop causal reasoning under complex scenarios. In contrast, Figure 5 reveals that the non-reasoning Qwen2.5-72B-Instruct over-generalises the fine-grained intent “Does this fall under student insurance?” into a broad liability query (predicting D instead of B), yet precisely recognises a straightforward capacity question, underscoring its sensitivity to subtle semantic boundaries and the need for tailored prompting or post-calibration to mitigate misclassification.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Business: 保险理解认知 / Insurance Understanding and Cognition(IUC)
Task: 多项选择 / Multiple choice

Prompt:
背景: 你是一名保险领域专家, 以下是针对保险条款进行属性抽取的单项选择题, 请直接给出正确答案的选项。
Background: As an expert in the insurance sector, you are presented with a multiple-choice question concerning attribute extraction from insurance policy terms. Kindly select the correct option.

保险条款:
重大疾病住院前后门急诊医疗费用 指被保险人经医院诊断罹患重大疾病必须接受住院治疗, 在住院前7日(含住院当日) 和出院后30日(含出院当日)内, 因与该次住院相同原因而接受重大疾病门急诊治疗时, 被保险人需个人支付的、必需且合理的治疗重大疾病门急诊医疗费用 (但不包括重大疾病特殊门诊医疗费用和重大疾病门诊手术医疗费用)。

Insurance Terms:
Major illness-related pre-hospitalization and post-hospitalization outpatient and emergency medical expenses pertain to costs that the insured must bear when diagnosed with a critical illness and hospitalized. These costs apply for treatments received within 7 days prior to (including the day of admission) and 30 days following (including the day of discharge) hospitalization for the same medical reason, excluding expenses associated with special outpatient care and outpatient surgical treatments for major illnesses.

抽取要求:
需要抽取的属性是: 【费用报销范围】; 该属性的类型是: 【ENUM】; 该属性的取值范围是: 【不限医保、医保外、医保内、医保内+医保外自费药、未知】

Extraction Requirements:
The attribute to be extracted is [Coverage Range]; type: [ENUM]; value range: [Not restricted by insurance, Outside insurance coverage, Within insurance coverage, Within insurance + self-paid medication outside insurance, Unknown]

选项:
A. 医保内+医保外自费药
B. 医保内
C. 医保外
D. 不限医保

Options:
A. Within insurance + self-paid medication outside insurance
B. Within insurance coverage
C. Outside insurance coverage
D. Not restricted by insurance

Answer:
D

Figure 6: A multiple-choice example.

Business: 保险安全合规 / Insurance Safety and Compliance (ISC)
Task: 判断 / True-False

Prompt
背景: 你是一名保险文案合规性审核专家, 你可以识别一段保险文案是否合规, 如果合规, 则回答“是”, 如果不合规, 则回答“否”。
Background: As a compliance review expert for insurance copy, you have the ability to determine its compliance status. Respond with 'Yes' if the copy is compliant, or 'No' if it is not.

不合规情况:
1. 出现歧视行为、特殊群体刻板印象词汇
2. 夸大产品保障
3. 宣传产品限时限量, 刺激购买
4. 引用不真实、不准确的数据和资料
5. 使用政府、监管等描述进行不当背书
6. 使用红包补贴、免费宣传进行诱导销售
7. 涉及法律功能不实描述
8. 使用绝对化用语、极限词
9. 隐瞒限制条件

Non-Compliance Assessment:
1. Discriminatory language or stereotypes about specific groups
2. Overstating the coverage provided by the product
3. Creating urgency by promoting time-limited or quantity-limited offers
4. Quoting inaccurate or misleading data and information
5. Inappropriately using endorsements related to government or regulatory authorities
6. Engaging in sales inducement through offers like discounts or free incentives
7. Misrepresenting legal functionalities
8. Using absolute or extreme language
9. Failing to disclose any limitations or conditions

保险文案:
随着生活节奏的加快, 很多人对于长期医疗保险有了新的认识和需求, 尤其是上班族, 面对日益增长的医疗费用, 如何加强自身的保障, 成为他们最关心的问题。好医保长期医疗(0免赔), 为你提供全方位的医疗保障, 每月保费11.77元起就可以享受最高400万的医疗保障, 日常小病小痛住院能报销, 重大疾病产生的住院医疗费用能报销, 社保外的进口药、自费药、靶向药以及120万的CAR-T细胞治疗都可以报销。不管你是年轻人还是老年人, 好医保长期医疗(0免赔)都可以为你提供全面的医疗保障。现在就来投保, 让你的生活更加安心!

Insurance Copy
As fast-paced living becomes the norm, there's a growing awareness and need for long-term medical insurance. Office workers, in particular, are increasingly concerned about rising medical costs and seek to enhance their coverage. The Good Health Insurance - Long-term Medical plan (with zero deductible) offers comprehensive medical protection. With premiums starting as low as 11.77 yuan per month, you can receive up to 4 million yuan in medical coverage. From minor illnesses to significant diseases, hospitalization expenses can be reimbursed. Coverage includes imported, self-funded, and targeted medications, as well as up to 1.2 million yuan for CAR-T cancer treatment not covered by social insurance. Whether young or old, Good Health Insurance - Long-term Medical (zero deductible) offers you the broadest medical coverage available. Sign up now for a more secure future!

Answer
否
False

Figure 7: A True/False example.

G EXAMPLES

To illustrate the heterogeneous nature of our evaluation framework, we showcase representative examples across three distinct task formats: a multiple-choice example as demonstrated in Figure 6, a True/False example as shown in Figure 7, a Q&A example as presented in Figure 8, and a multi-turn dialogue example as presented in Figure 9.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



Figure 8: A single-turn Q&A example.



Figure 9: A multi-turn dialogue example.

	数据获取 /Data Collection
1026	
1027	
1028	
1029	## 任务: 将以下提供的保险领域文本内容, 转化为结构化的问答 (Q&A) 形式。
1030	(知识点)
1031	...
1032	## 要求:
1033	1. **提取关键信息:** 准确识别并提取文本中关于保险概念、定义、流程、特点、条款或案例等的关键知识点。
1034	2. **转化为问答对:** 将每个关键知识点设计为一个清晰的问题和相应的答案。
1035	3. **确保信息准确性:** 问答内容必须忠实于原文, 不得增加、删除或曲解原文提供的信息。
1036	4. **问题明确易懂:** 设计的问题应直接指向某个知识点, 语言简洁, 易于非专业人士理解 (如果适用)。
1037	5. **答案简洁明了:** 答案应直接回答问题, 内容精炼, 提取原文的核心解释。
1038	6. **覆盖主要内容:** 确保文本中的主要知识点都得到了覆盖, 形成相应的问答。
1039	7. **格式规范:** 以清晰的问答列表形式呈现, 例如使用 "问题:" 和 "答案:" 作为每一对问答的标记。
1040	## 输出规范: 以json格式输出: {"question": "问题", "answer": "答案"}
1041	## 待处理文本:
1042	...
1043	{text}
1044	...
1045	
1046	
1047	
1048	
1049	
1050	
1051	
1052	
1053	
1054	
1055	
1056	
1057	
1058	
1059	
1060	
1061	
1062	
1063	
1064	
1065	
1066	
1067	
1068	
1069	
1070	
1071	
1072	
1073	
1074	
1075	
1076	
1077	
1078	
1079	

Figure 10: The prompt for converting knowledge to questions.

	数据增强 /Data Augmentation
1049	
1050	## 任务1: 对以下文本进行同义词改写。
1051	## 要求:
1052	1.保持原意不变: 改写后的文本必须准确传达原文的核心意义和信息。
1053	2.使用同义词或近义词: 替换原文中的词语和短语, 使用意义相近但不同的词汇。
1054	3.改变表达方式: 可以在词语替换的基础上, 适当调整句子结构或表达方式, 使文本更具多样性。
1055	4.确保流畅自然: 改写后的文本应语法正确, 句子连贯, 阅读起来自然流畅, 避免生硬的词语堆砌。
1056	5.考虑上下文: 选择同义词时需考虑其在上下文中的适用性。
1057	## 待改写文本:
1058	...
1059	{text}
1060	...
1061	
1062	
1063	
1064	
1065	
1066	
1067	
1068	
1069	
1070	
1071	
1072	
1073	
1074	
1075	
1076	
1077	
1078	
1079	

Figure 11: The prompt for data augmentation.

H PROMPTS

H.1 PROMPTS FOR DATA CONSTRUCTION

Figure 10 shows the prompt for converting knowledge to questions. Figure 11 shows the prompt for data augmentation. Figure 12 shows the prompt for data quality control.

H.2 PROMPTS FOR EVALUATION

Figure 13 and Figure 14 shows the prompt for answer content splitting and matching for calculating faithfulness.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133



Figure 12: The prompt for data quality control.



Figure 13: The prompt for answer content splitting.

```

1134 ## 任务描述
1135 你的任务是根据给定的context判断statements中所有陈述句的忠诚度。对于每个陈述句，如果该陈述句可以直接基于context推理出来或者检索出来，则返回判定为1；如果该陈述句不能直接基于context推理或者检索出来，则返回判定为0。请直接输出结果，除此以外不要输出其它内容。
1136 ## Task Description
1137 Your task is to assess the fidelity of statements based on the given context. For each statement, assign a value of 1 if it can be directly inferred or sourced from the context, or 0 if it cannot.
1138 Please output only the results without any additional content.
1139 输入:
1140 Input:
1141 {
1142   "context": [text],
1143   "statements": [text]
1144 }
1145 ## 参考案例
1146 输入:
1147 {
1148   "context": ["检索1线上投保通常可以选择多种便捷的付款方式，包括银行卡支付（借记卡、信用卡）、第三方支付平台（如支付宝、微信支付）、以及保险公司的官方APP或网站绑定的支付渠道。一些保险公司还支持分期支付以减轻一次性付款压力。选择付款方式时要确保符合保险公司要求，并注意保费扣款时间，以避免因支付失败导致保单生效受影响。此外，还请注意保存付款凭证以备后续查询或理赔使用。"],
1149   "statements": ["检索2: 保费可采取多种支付方式，通常包括现金支付、银行卡转账支付、支付宝支付等。", "检索3: 确保信用卡在可用方式列表中。", "检索4: 综上，信用卡支付是可行的。", "检索5: 但需结合具体投保渠道的设置。"]
1150 }
1151 输出:
1152 [{"statement": "确保信用卡在可用方式列表中。", "verdict": 1},
1153 {"statement": "综上，信用卡支付是可行的。", "verdict": 1},
1154 {"statement": "但需结合具体投保渠道的设置。", "verdict": 1}]
1155 ## Reference Example
1156 Input:
1157 {
1158   "context": ["Source 1: Online insurance enrollment typically allows for various easy payment options, such as bank card payments (both debit and credit cards), third-party platforms (like Alipay, WeChat Pay), and payment channels linked to the insurance company's app or website. Installment plans are also available with some insurers to alleviate the pressure of full upfront payment. Ensure compliance with insurer payment guidelines and be mindful of payment schedules to prevent issues with policy activation due to payment hiccups. It's also important to keep the payment receipt handy for future inquiries or claims.\nSource 2: Premiums can be settled through multiple payment methods, commonly including cash, bank card transfers, and Alipay. ",
1159   "statements": ["Make sure that credit card payment is an option listed among the available methods.", "In conclusion, using a credit card for payment is possible.", "Nonetheless, it should align with the specific arrangements of the insurance purchase platform."]
1160 }
1161 Output:
1162 [{"statement": "Make sure that credit card payment is an option listed among the available methods.", "verdict": 1},
1163 {"statement": "In conclusion, using a credit card for payment is possible.", "verdict": 1},
1164 {"statement": "Nonetheless, it should align with the specific arrangements of the insurance purchase platform.", "verdict": 1}]

```

Figure 14: The prompt for answer matching.

```

1158 ## 任务描述
1159 给你一个question和answer，请结合question信息，并分析answer文本的复杂度，将answer的文本分解成最多20个可以完全可理解的陈述句，确保任何陈述句中不使用代词。陈述句的内容必须是answer的原文内容，不得进行任何隐含的数学计算或推导。请直接输出最终结果，除此以外不要输出其它内容。
1160 输入:
1161 {
1162   "question": "0",
1163   "answer": {}
1164 }
1165 ## 参考案例
1166 输入:
1167 {
1168   "question": "平安中老年人意外险特定区域和其他区域的核心差异有哪些？",
1169   "answer": "您好，平安中老年人意外险的特定区域包括公共交通工具内部和公益性文化设施内，其他区域则是这些地方之外的普通区域。在特定区域内发生意外伤害，不同套餐的保额会更高，比如升级版在特定区域身故残疾保额能达到20万，其他区域是10万。另外，意外医疗报销的免赔额是100元，升级后免赔额是0元。升级后在特定区域身故残疾保额为10万，意外医疗报销的免赔额是100元。升级后，意外医疗按80%赔付。", "未投保报销时，意外医疗按60%赔付。", "意外医疗报销规则在特定区域和其他区域都适用。"]
1170 }
1171 输出:
1172 [{"statement": "The specific areas of Ping An Middle-aged and Elderly Accident Insurance include the interiors of public transport vehicles and public welfare cultural facilities.", "The other areas of Ping An Middle-aged and Elderly Accident Insurance are ordinary areas outside the specific areas.", "When an accident occurs in a specific area, the insurance amount of different packages will be higher.", "The death and disability insurance amount of the upgraded package can reach 200,000 yuan in specific areas.", "The death and disability insurance amount of the upgraded package is 100,000 yuan in other areas.", "The deductible for accident medical reimbursement is 100 yuan.", "After social security reimbursement, accident medical expenses are compensated at 80%.", "When not reimbursed through social security, accident medical expenses are compensated at 60%. This applies to both types of areas."}]
1173 [{"statement": "The specific areas of Ping An Middle-aged and Elderly Accident Insurance include the interiors of public transport vehicles and public welfare cultural facilities.", "The other areas of Ping An Middle-aged and Elderly Accident Insurance are ordinary areas outside the specific areas.", "When an accident occurs in a specific area, the insurance amount of different packages will be higher.", "The death and disability insurance amount of the upgraded package can reach 200,000 yuan in specific areas.", "The death and disability insurance amount of the upgraded package is 100,000 yuan in other areas.", "The deductible for accident medical reimbursement is 100 yuan.", "After social security reimbursement, accident medical expenses are compensated at 80%.", "When not reimbursed through social security, accident medical expenses are compensated at 60%. This applies to both types of areas."}]

```

Figure 15: The prompt for ground truth content splitting

Figure 15 and Figure 16 shows the prompt for ground truth content splitting and matching for calculating completeness.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

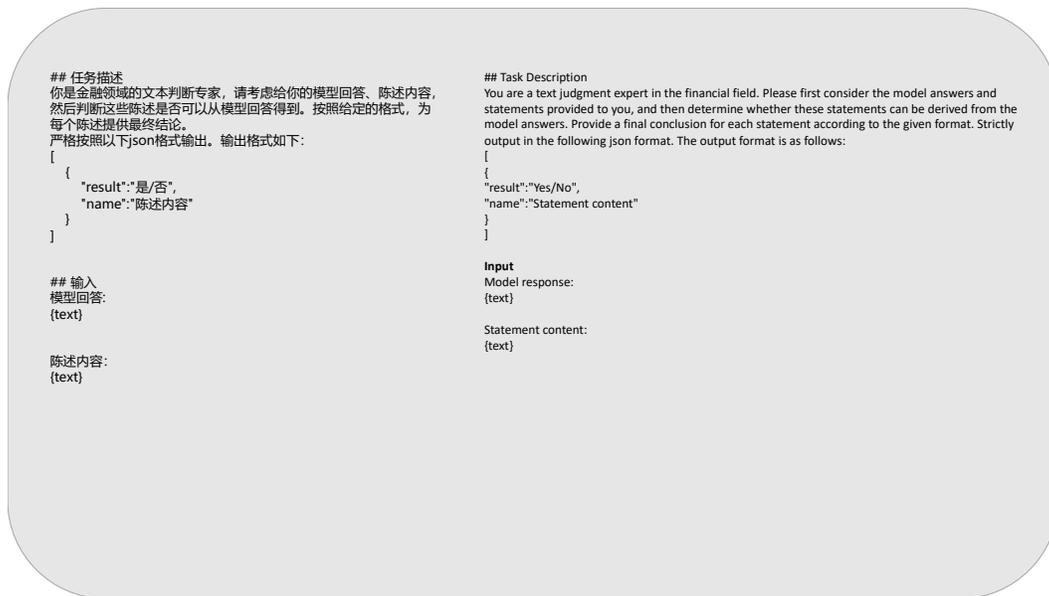


Figure 16: The prompt for ground truth matching.