

LINGLANMiDIAN: SYSTEMATIC EVALUATION OF LLMs ON TCM KNOWLEDGE AND CLINICAL REASONING

Rui Hua^{†,1}, Yu Wei^{†,2,3,4}, Zixin Shu^{†,5,6,7}, Kai Chang¹, Dengying Yan¹, Jianan Xia¹, Zeyu Liu¹
 Hui Zhu^{5,6,7}, Shujie Song^{5,6,7}, Mingzhong Xiao^{5,6,7}, Xiaodong Li^{5,6,7}, Dongmei Jia¹⁵, Zhuye Gao¹⁵, Yanyan Meng¹¹
 Naixuan Zhao¹², Yu Fu¹³, Haibin Yu¹⁴, Benman Yu¹⁴, Yuanyuan Chen¹⁴, Fei Dong¹⁷, Zhizhou Meng¹³
 Pengcheng Yang^{2,3,4}, Songxue Zhao^{2,3,4}, Lijuan Pei^{2,3,4}, Yunhui Hu^{2,3,4}, Kan Ding¹⁰, Jiayuan Duan¹⁸
 Wenmao Yin¹⁹, Yang Gu¹⁶, Runshun Zhang⁸, Qiang Zhu¹, Jian Yu¹, Jiansheng Li¹⁴
 Baoyan Liu⁹, Wenjia Wang^{*,2,3,4}, Xuezhong Zhou^{*,1}

¹Department of Computer Science and Technology, Beijing Jiaotong University, Beijing, China

²Tianjin Tasly Digital Chinese Medicine Technology Co., Ltd., Tianjin, China

³Tasly Biopharmaceuticals Co., Ltd., Tianjin, China

⁴State Key Laboratory of Chinese Medicine Modernization, Tianjin, China

⁵Institute of Liver Diseases, Hubei Key Laboratory of the theory and application research of liver and kidney in traditional Chinese medicine, Hubei Provincial Hospital of Traditional Chinese Medicine, Wuhan, China

⁶Affiliated Hospital of Hubei University of Chinese Medicine, Wuhan, China

⁷Hubei Province Academy of Traditional Chinese Medicine, Wuhan, China

⁸Department of Gastroenterology, Guang'anmen Hospital, China Academy of Chinese Medical Sciences, Beijing, China

⁹China Academy of Chinese Medical Sciences, Beijing, China

¹⁰China Institute for History of Medicine and Medical Literature, China Academy of Chinese Medical Sciences, Beijing, China

¹¹Beijing Research Institute of Chinese Medicine, Beijing University of Chinese Medicine, Beijing, China

¹²Beijing University of Chinese Medicine Third Affiliated Hospital, Beijing University of Chinese Medicine, Beijing 100029, China

¹³School of Traditional Chinese Medicine, Beijing University of Chinese Medicine, Beijing, China

¹⁴The First Affiliated Hospital, Henan University of Chinese Medicine, Zhengzhou, China

¹⁵Xiyuan Hospital, China Academy of Chinese Medical Sciences, Beijing, China

¹⁶Da'an Health Technology (Beijing) Co., Ltd, Beijing, China

¹⁷Institute of Chinese Medicine Epidemic Disease, Beijing University of Chinese Medicine, Beijing, China

¹⁸Beijing Zhongtengbaimai Medical Technology Co., Ltd, Beijing, China

¹⁹Sinsoft Company Limited, Beijing, China

[†] These authors contributed equally to this work.

* Corresponding author.

ABSTRACT

Large language models (LLMs) are advancing rapidly in medical NLP, yet Traditional Chinese Medicine (TCM) with its distinctive ontology, terminology, and reasoning patterns requires domain-faithful evaluation. Existing TCM benchmarks are fragmented in coverage and scale and rely on non-unified or generation-heavy scoring that hinders fair comparison. We present the LingLanMiDian (LingLan) benchmark, a large-scale, expert-curated, multi-task suite that unifies evaluation across knowledge recall, multi-hop reasoning, information extraction, and real-world clinical decision-making. LingLan introduces a consistent metric design, a synonym-tolerant protocol for clinical labels, a per-dataset 400-item Hard subset, and a reframing of diagnosis and treatment recommendation into single-choice decision recognition. We conduct comprehensive, zero-shot evaluations on 14 leading

open-source and proprietary LLMs, providing a unified perspective on their strengths and limitations in TCM commonsense knowledge understanding, reasoning, and clinical decision support; critically, the evaluation on Hard subset reveals a substantial gap between current models and human experts in TCM-specialized reasoning. By bridging fundamental knowledge and applied reasoning through standardized evaluation, LingLan establishes a unified, quantitative, and extensible foundation for advancing TCM LLMs and domain-specific medical AI research. All evaluation data and code are available at <https://github.com/TCMAI-BJTU/LingLan> and <http://tcmlnp.com>.

Keywords Traditional Chinese Medicine, Large Language Models, Benchmark, Evaluation

1 Introduction

Foundation models trained on web scale corpora have redefined language understanding and reasoning, with rapid advances across capability and scale [1, 2, 3, 4]. This shift from task specific systems to general purpose language intelligence has motivated unified evaluations that prioritize breadth and comparability, exemplified by MMLU and HELM [5, 6]. Contemporary model families such as the GPT-4 and GPT-5 series, DeepSeek, and Qwen illustrate a trajectory toward stronger multistep reasoning and cross domain transfer [2, 7, 8].

In medicine, frontier LLMs can approach or even surpass human performance on standardized knowledge assessments such as MedQA, CMB, and CMExam [9, 10, 11], alongside progress in biomedical and clinical modeling (for example, BioGPT and Med-PaLM) and applied evaluation of clinical use cases [12, 13, 14, 15]. Clinical reasoning, however, demands more than factual recall and requires the integration of structured knowledge with contextual patient data and accountable decision processes [16, 17, 18]. Recent surveys therefore emphasize multi-dimensional and clinically grounded evaluation beyond narrow QA, while noting that most benchmarks remain anchored in modern biomedicine [19]. Within the Chinese medical ecosystem, general domain and biomedical LLMs such as Baichuan, Huatuo, and Zhongjing have been adapted to consultation and decision support [20, 21, 22]. Building on these efforts, TCM-oriented models such as TCMChat and Lingdan fine tune on curated corpora spanning classical canons, syndrome and treatment guidelines, and prescription data to better reflect Traditional Chinese Medicine’s knowledge system and reasoning paradigm [23, 24].

Across the history of AI, benchmark datasets have played a pivotal role in catalyzing progress. ImageNet revolutionized computer vision by establishing a large-scale, standardized benchmark that defined measurable progress and spurred model innovation [25]. Analogously, language benchmarks such as MMLU [5], C-Eval [26], and HELM [6] have become cornerstones for evaluating reasoning, fairness, and calibration in LLMs. In the medical domain, datasets such as MedQA [9], PubMedQA [27], and CMB [10] have enabled rigorous, reproducible assessment of medical LLMs. However, no existing benchmark provides equivalent comprehensiveness for TCM, which represents a culturally rooted and conceptually distinct medical system, thereby limiting the systematic advancement of AI models in this domain.

TCM embodies a deeply empirical and experience-driven medical paradigm. Unlike the protocolized frameworks of Western biomedicine, TCM emphasizes the recognition of dynamic patterns rather than fixed disease entities, and the corresponding therapeutic principle is determined through interpretive reasoning grounded in accumulated clinical experience rather than strictly codified rules [28, 29]. This experience-centered reasoning extends to prescription formulation, where herbal combinations are adjusted according to subtle patient variations, temporal changes, and practitioner intuition. Such contextual flexibility makes TCM both highly adaptive and inherently difficult to formalize. From a computational perspective, these characteristics create unique challenges for model evaluation: the same clinical presentation may correspond to multiple legitimate syndromes or treatment strategies, while similar herbal prescriptions can differ in dosage proportions yet remain therapeutically equivalent. Moreover, TCM texts often employ metaphorical and non-literal expressions rooted in classical Chinese, which complicate token-level language modeling and semantic normalization. As a result, an effective benchmark must assess not only factual knowledge retrieval but also models’ ability to approximate the interpretive, experience-based reasoning that underpins authentic TCM practice.

Despite growing efforts to build TCM-oriented datasets, most existing benchmarks capture only a narrow slice of this experiential reasoning process. Resources such as TCMBench [30], TCMD [31], and MTCMB [32] predominantly focus on examination-style or text-generation tasks, emphasizing factual recall rather than the process of clinical synthesis. Furthermore, most corpora are derived from publicly available texts or restructured educational materials rather than directly curated clinical sources, resulting in limited representation of nuanced practitioner experience. In addition, existing benchmarks often assume deterministic correctness, applying strict exact-match metrics to problems that are inherently probabilistic and context-sensitive. This methodological gap neglects the experiential uncertainty fundamental to TCM practice, where flexible equivalence, encompassing synonymous herbs, proportional dosage variation, or overlapping syndromes, is clinically acceptable. Consequently, current benchmarks fail to reflect the interpretive and adaptive reasoning that defines expert-level TCM decision-making. A comprehensive evaluation

framework must therefore integrate structured and experiential components, enabling fair comparison of LLMs not only by recall accuracy but also by their ability to simulate human-like diagnostic and therapeutic judgment.

We introduce LingLanMiDian (LingLan), a rigorously curated benchmark that elevates evaluation of LLMs for TCM from factual recall to clinically salient reasoning. LingLan provides broad domain coverage, expert-audited content, hard subsets for robustness, and a unified metric suite that enables commensurate, reproducible comparison across heterogeneous task formats. Under a standardized zero-shot protocol, results show near-ceiling performance on licensing-style recall but persistent deficits in multi-hop synthesis, prescription composition, and dosage proportionality, underscoring the gap between surface knowledge and expert-level TCM reasoning.

Our contributions are as follows:

- (1) A comprehensive, expert-curated TCM suite spanning five domains and 13 subtasks.
- (2) A unified, reproducible metric framework with task-aligned scoring and difficulty-calibrated hard subsets.
- (3) Systematic cross-model baselines that expose actionable failure modes in option calibration, information extraction scaling, and clinical reasoning.

2 Related Work

To contextualize LingLan, we survey prior TCM and clinical NLP benchmarks and identify enduring gaps in dataset scale, task breadth, and harmonized metrics that our benchmark is intended to address.

2.1 General Domain Benchmarks

In the general domain, several large-scale benchmarks probe knowledge and reasoning under controlled conditions. C-Eval [26] assembles 13,948 multiple-choice questions spanning 52 disciplines and four difficulty tiers (middle school, high school, college, professional), accompanied by a “hard” split and accuracy-based evaluation, providing a broad, curriculum-aligned stress test for Chinese. GPQA [33] targets graduate-level, “Google-proof” questions crafted to resist retrieval and superficial pattern matching, thereby emphasizing depth of domain expertise and non-trivial reasoning. Humanity’s Last Exam [34] proposes an adversarial, web-resistant evaluation paradigm that couples difficult expert-authored items with rigorous auditing to assess reasoning quality and safety beyond conventional QA. MATH [35] compiles competition-style mathematics problems with step-by-step solutions and shows that accuracy remains low even as model scale increases, underscoring the limits of brute-force scaling for formal reasoning. Collectively, these resources motivate comprehensive, difficulty-calibrated, and retrieval-resistant evaluation—principles that inform the design of LingLan for TCM.

2.2 Modern Biomedical Benchmarks

Beyond general domains, a number of general medical benchmarks have been proposed to assess the medical reasoning capabilities of LLMs. CMExam [11] comprises 60K+ multiple-choice questions drawn from the Chinese National Medical Licensing Examination, accompanied by official-style solution explanations and five expert-labeled attributes (disease group, department, discipline, competency area, difficulty), and is typically evaluated by accuracy. MedQA [9] collects professional board-exam questions across three languages, with 12,723 English, 34,251 Simplified Chinese, and 14,123 Traditional Chinese examples, forming a multiple-choice OpenQA resource widely used for accuracy-based evaluation. CMB [10] offers a comprehensive Chinese medical benchmark combining 280,839 multiple-choice questions (CMB-Exam) and 74 expert-curated clinical consultation cases (CMB-Clin), compiled from national licensing exams, textbooks, and teaching materials with standardized evaluation splits. MedBench [36] provides 40,041 questions sourced from authentic Chinese medical examinations (e.g., licensing, resident training, doctor-in-charge) and real clinical reports, supporting unified accuracy-style assessments for knowledge and diagnostic reasoning. Palepu et al. [37] introduce the “Humanity’s Next Medical Exam” framework, which redefines medical AI evaluation by comprising three pillars, including interactive interrogation, sandboxed experiential learning, and real-world continuous learning, and argue that QA-centric benchmarks fail to capture clinical complexity. MedHELM [38] introduces a clinician-validated, practice-grounded evaluation suite (35 benchmarks over 121 tasks) for nine state-of-the-art models, showing reasoning-centric models often lead but task- and cost-dependent trade-offs persist.

However, these resources are largely grounded in modern biomedicine and general clinical paradigms, offering limited coverage of TCM’s unique ontology, terminology, and reasoning frameworks. LingLan differentiates itself by providing a TCM-specific benchmark that unifies standardized examination-style questions, foundational knowledge, and applied clinical reasoning grounded in both ancient textual sources and modern clinical practice.

Table 1: Comparison of major TCM benchmarks by scale, task types, sources, and metrics. Abbreviations: SC (single-choice), MC (multiple-choice), Cloze (fill-in-the-blank), IE (information extraction), DTR (diagnostic–therapeutic reasoning), DR (decision recognition), TLE (TCM Licensing Examination), FTK (Fundamental TCM Knowledge), CPMK(Chinese Patent Medicine Knowledge), EMR(electronic medical records), Char-F1(character-level F1), MAE(mean absolute error).

Benchmark	Scale	Question Type	Sources	Metrics
TCMD [31]	3,451	MC	Licensing examinations	Accuracy
TCMBench [30]	5,473	MC	Licensing examinations	TCMScore, ROUGE, BERTScore
MTCMB [32]	7,100 (6,000 exam samples)	MC / Cloze / Open QA / Prescription / Diagnosis / IE	Licensing examinations / Classical texts / Clinical records / Guidelines	Accuracy, ROUGE, BLEU, BERTScore
TCMEval-SDT [39]	300 cases	IE / Pathogenesis / Syndrome / Summary	Internet / Classical texts / Clinical records	Stage-wise qualitative
LingLan (ours)	25,624 (~2,000 per subtask)	SC / MC / Cloze / IE / DTR / DR	TLE / FTK / CPMK / EMR / Classical texts / Master-physician casebooks	Accuracy, Precision, Recall, Char-F1, MAE, Cosine similarity

2.3 Traditional Chinese Medicine Benchmarks

Efforts to evaluate LLMs within the domain of TCM have only emerged in recent years, highlighting the increasing demand for domain-specific and standardized evaluation frameworks. Early studies mainly focused on examination-style question answering. Currently, as shown in Table 1, these benchmarks provide an initial comparison across models, but still leave gaps in broader, more comprehensive evaluation.

TCMD [31] compiles 3,451 CNMLE-style multiple-choice questions (2,851 train / 600 test) with accompanying explanations; items are OCR-transcribed, de-duplicated, and human-verified under the official exam manual to balance subjects, and evaluation is reported with accuracy.

TCMBench [30] centers on the TCM-ED question set built from the TCM Licensing Examination, reporting around 5,473 QA pairs in the paper and an updated 6,482 QA pairs in the official repository, of which 1,300 include official standard explanations; questions were expert-filtered and student-checked, and answers are scored by TCMScore alongside ROUGE/BERTScore.

MTCMB [32] co-develops 12 sub-datasets across five categories: knowledge QA, language understanding, diagnostic reasoning, prescription recommendation, and safety. The benchmark totals around 7,100 examples; sources span licensing exams, classical texts, and real or simulated clinical records, with expert curation and a mix of accuracy and generation-based metrics (e.g., ROUGE, BLEU, and BERTScore).

TCMEval-SDT [39] provides a curated set of 300 expert-annotated cases for syndrome differentiation, sourced from internet repositories, classical TCM literature, and de-identified hospital records with FAIR-compliant metadata. The benchmark formalizes a four-stage evaluation protocol comprising information extraction, pathogenesis inference, syndrome reasoning, and explanatory summarization.

LingLan complements these benchmarks by integrating a much larger sample size (over 25,000 items), covering both structured and unstructured text scenarios, and offering standardized, quantitative evaluation across domains.

2.4 LLMs for Chinese Medical and TCM Applications

Huatuo [21] fine tunes LLaMA-7B [3] with Chinese medical knowledge grounded QA generated from CMeKG to improve reliability on biomedical queries and introduces an evaluation metric for safety, usability, and smoothness. Zhongjing [22] adopts an expert in the loop framework that learns from real world multi turn Chinese medical dialogues to align models with clinical reasoning and safety, achieving leading performance on Chinese medical QA and dialogue benchmarks. Lingdan [24] enhances the encoding of TCM knowledge for clinical reasoning tasks using LLMs, improving diagnosis, treatment, and prescription recommendation. TCMChat [23] builds on Baichuan2-7B-Chat [20] and is trained with pretraining and supervised fine tuning on curated TCM corpora and Chinese QA datasets spanning six task types. BianCang [40] develops a TCM oriented model through full parameter fine-tuning, which involves continuing pretraining and supervised alignment on Qwen2.5-7B using hospital EMRs and knowledge derived from the Chinese Pharmacopoeia, and reports gains across diverse TCM benchmarks. Despite these advances, most other Chinese medical and TCM systems rely on backbones below 14B parameters and use adapter based methods such as

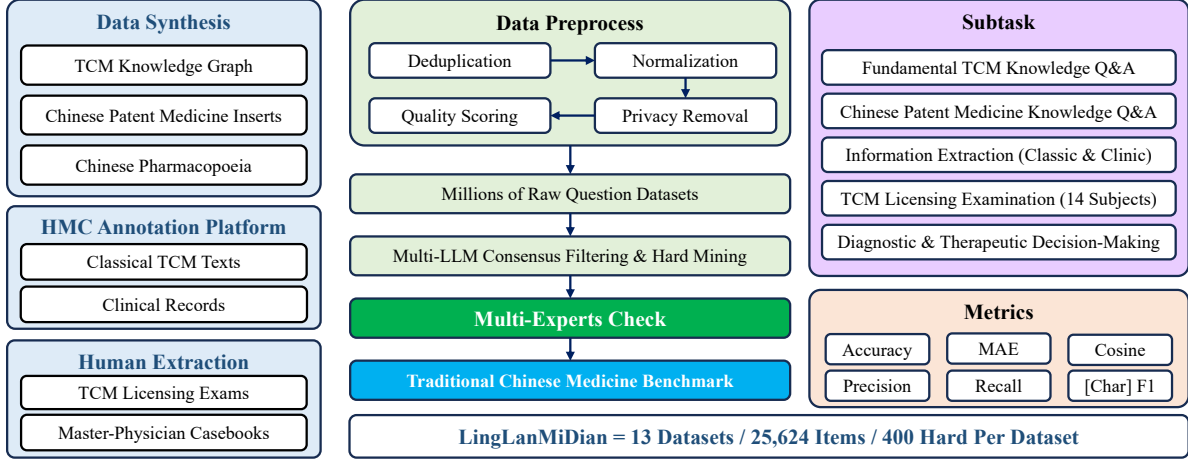


Figure 1: Overview of the LingLan construction pipeline. The left panel summarizes data sources; the center panel illustrates data processing and curation; the right panel presents the task taxonomy covered in LingLan together with the unified evaluation metrics. Abbreviations: HMC (Human-machine collaborate).

LoRA rather than full parameter fine tuning, and their absolute performance remains clearly below that of contemporary large capacity models such as the DeepSeek family.

Existing TCM benchmarks have made important progress toward evaluating LLMs in this domain, but they suffer from three major limitations: (1) insufficient task diversity, with most focusing on question answering or generative tasks; (2) reliance on subjective scoring methods, limiting reproducibility; and (3) small scale and lack of standardized, fine-grained metrics. LingLan addresses these gaps by offering a large-scale, multi-task, and expert-curated evaluation framework that unifies objective metrics across both structured and reasoning-intensive tasks, establishes an expert-reviewed hard subset for challenge evaluation, and provides the first holistic assessment of TCM reasoning across 14 state-of-the-art LLMs.

3 Methodology

We now detail LingLan’s construction, covering its data sources and curation, task taxonomy, metric definitions, and standardized zero-shot evaluation settings, to provide a reproducible blueprint for future extensions.

3.1 Overview

LingLanMiDian (LingLan) benchmark is designed to systematically evaluate the performance of LLMs across the full spectrum of TCM knowledge and reasoning. It unifies heterogeneous evaluation tasks such as knowledge recall, multi-hop reasoning, information extraction, and diagnostic decision-making under a consistent and quantitative framework, as illustrated in Figure 1. LingLan comprises 13 datasets spanning five domains: *Licensing Examination*, *Fundamental TCM Knowledge*, *Chinese Patent Medicine*, *Information Extraction*, and *Diagnostic and Therapeutic Decision-Making*. In total, the benchmark contains 25,624 expert-verified items, with approximately 2,000 samples per dataset and an additional 400-item *Hard subset* in each domain to assess advanced reasoning capability.

3.2 Data Sources

(1) TCM Licensing Examination (TLE). Questions were collected from the official National TCM Qualification Examination (2000–2016), covering a comprehensive range of medical subjects including *Fundamentals of TCM Theory*, *Diagnostics of TCM*, *Chinese Materia Medica*, *Formulary Science*, *Internal Medicine of TCM*, *Surgery of TCM*, *Gynecology of TCM*, *Pediatrics of TCM*, and *Acupuncture and Moxibustion*. In addition, the dataset integrates supporting courses required in medical licensing education, such as *Fundamentals of Diagnostics*, *Internal Medicine*, *Infectious Diseases*, *Medical Ethics*, and *Health Regulations*. Together, these subjects comprehensively represent the standardized knowledge framework assessed in national medical education and form the foundation for evaluating both factual and procedural competence in TCM.

(2) Fundamental TCM Knowledge (FTK). This dataset is grounded in a self-constructed TCM knowledge graph assembled from authoritative sources, including canonical classical texts, contemporary peer-reviewed TCM/biomedical literature, and nationally adopted teaching materials, and further augmented with public databases such as SymMap [41] and TCMID [42]. The graph consolidates and normalizes entities and relations across herbs, ingredients, efficacies, symptoms, syndromes, prescriptions, symptom clusters, diseases, and genes, capturing hierarchical and therapeutic links (e.g., herb–ingredient, herb–efficacy, symptom–syndrome, syndrome–prescription) under a unified ontology.

(3) Chinese Patent Medicine Knowledge (CPMK). This dataset is built from structured knowledge extracted from the Pharmacopoeia of the People’s Republic of China and package inserts of Chinese patent medicines. The source corpus was normalized into a unified schema covering herbal composition, associated diseases and syndromes, pharmacology and toxicology, indications, contraindications, product specifications, safety information, etc. In total, it comprises standardized records for more than 9,000 Chinese patent medicines.

(4) Information Extraction (IE). The information extraction datasets draw from two sources: classical TCM literature written in Literary Chinese and real-world clinical medical records composed in modern vernacular Chinese. Annotated samples were produced via a semi-automated pipeline that combines named entity recognition (NER) algorithms with expert verification [43]. All clinical texts were fully de-identified prior to annotation to ensure privacy compliance. The resulting corpora encompass more than one hundred entity categories, including positive and negative symptoms, herbs, positive and negative physical signs, condition changes, and episode characteristics. This dual-register design captures both historical and contemporary linguistic styles, enabling fine-grained evaluation of LLMs in structured information extraction and clinical text understanding across both classical and modern contexts.

(5) Diagnostic and Therapeutic Decision-Making. The diagnostic and therapeutic datasets are built from more than 13,000 clinical case records authored by over 300 renowned TCM physicians, covering 7,595 syndromes. Domain experts extracted and verified high-quality annotations for syndrome differentiation, treatment principles, herbal prescriptions, and corresponding dosages. We further assessed coverage against the ICD-11 [44] “Traditional medicine patterns” category: excluding “specified” and “unspecified” terms, there are 216 categories in total, of which the dataset covers 174, yielding a coverage rate of 80.56%. These resources enable holistic evaluation of LLMs on integrated clinical reasoning and decision-making in TCM.

Except for TLE, all datasets are non-public or controlled constructions that do not appear in internet corpora and, in principle, have not been observed during pretraining by existing LLMs.

FTK		CPMK		TLE
Single-choice 1,962	Cloze 1,959	Single-choice 2,000	Cloze 2,000	Single-choice 1,832
Multi-choice 1,919		Multi-choice 2,000		
DR (SC)		IE		DTR
Syndrome Diagnosis 2,000	Prescription 2,000	Classical Texts 2,000	Clinical Texts 2,000	Generation 2,000
Treatment Method 2,000				

Figure 2: Subtask distribution and dataset sizes in LingLan. Abbreviations: FTK (Fundamental TCM Knowledge), CPMK (Chinese Patent Medicine Knowledge), TLE (TCM Licensing Examination), DR (Decision Recognition), IE (Information Extraction), DTR (Diagnostic–Therapeutic Reasoning), SC (single-choice).

3.3 Data Curation and Quality Assurance

Across all domains, LingLan applied a unified preprocessing pipeline including text normalization, de-duplication, automatic quality screening, and expert verification. All clinical materials were de-identified to ensure privacy compliance, and low-quality or ambiguous entries were systematically excluded prior to inclusion.

3.3.1 Task Definition and Taxonomy

LingLan establishes a unified taxonomy encompassing 13 subtasks (as shown in Figure 2) across five major domains of TCM, with a total of 44 quantitative evaluation dimensions. Each subtask is designed to capture distinct dimensions of

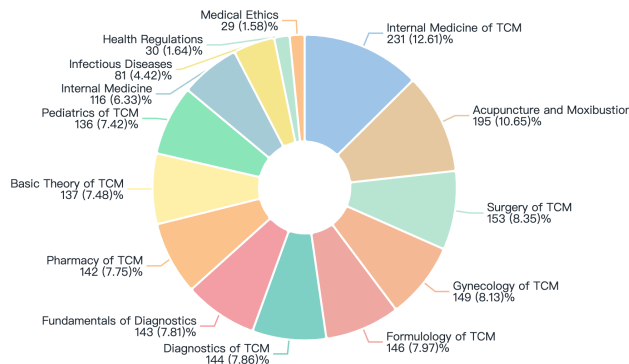


Figure 3: Distribution of items across the 14 subjects in the TLE (TCM Licensing Examination) dataset.

model capability, ranging from factual recall and symbolic association to structured extraction and clinical reasoning, thereby ensuring a balanced and multi-perspective assessment framework. All subtasks adopt standardized input-output formats and metric definitions, enabling cross-model comparability and reproducible benchmarking.

3.3.2 Expert Annotation and Verification

All datasets in LingLan underwent a rigorous expert validation process to ensure factual accuracy, terminological precision, and consistency with established principles of TCM. The review framework was designed to eliminate ambiguity, enforce theoretical correctness, and guarantee linguistic clarity. During annotation and verification, experts were prohibited from consulting LLMs or other automated assistants to avoid bias. Instead, they consulted authoritative references, including the Pharmacopoeia of the People’s Republic of China, Differential Diagnosis of TCM Symptoms, and national TCM Teaching Materials. The review guidelines emphasized semantic precision and task-specific rigor. To ensure question clarity, all stems were revised to remove ambiguity and to align precisely with the intended evaluation target; for example, vague fill-in-the-blank items were reformulated to specify whether the expected answer is a syndrome or a symptom. Terminological accuracy was enforced by correcting misuse of TCM terms (such as conflating syndromes with diseases or symptoms) to canonical phrasing. For answer validation, reference answers were verified against standard textbooks and authoritative sources; when multiple answers were acceptable, they were listed with enumeration marks to preserve complete semantic coverage. For option revision, ambiguous or partially correct multiple-choice options were replaced with unambiguously correct or incorrect alternatives to enhance discriminative validity. As part of quality control, items with irreparable logical or semantic flaws were excluded from the final dataset. This standardized protocol ensures that LingLan items are theoretically sound, linguistically precise, and pedagogically robust, thereby minimizing uncertainty and maximizing reproducibility across diverse LLM evaluation settings.

3.3.3 Knowledge-Oriented Tasks

This subsection describes the curation and quality assurance for the three knowledge-oriented suites, including the TCM Licensing Examination (TLE), Fundamental TCM Knowledge (FTK), and Chinese Patent Medicine Knowledge (CPMK).

For TLE, all items were standardized into a single-choice format with a unique gold answer and a normalized subject tag. We removed duplicates through fingerprinting of stems and explanations, corrected typography and punctuation, and rewrote unclear stems to eliminate answer ambiguity. Option sets were audited to ensure one and only one valid key, and distractors that were partially correct or definition-overlapping were replaced. Items that could not be unambiguously repaired were discarded, yielding a final pool of 1,832 high-quality questions. As shown in Figure 3, the corpus spans 14 subjects.

For the FTK task, we synthesize questions from an in-house TCM knowledge graph by systematically traversing its triples. For each head entity, we retrieve its associated triples in fixed-size batches of 20, when an entity is linked to more than 20 triples, we iterate over successive, non-overlapping batches of size 20, treating each batch as an independent background context. Using task-specific prompts and the Qwen3-32B model, we generate a large pool of candidate items spanning single-choice, multiple-choice, and cloze (fill-in-the-blank) formats, totaling several hundred thousand entries. For the CPMK task, we start from a curated knowledge base of Chinese patent medicines and use, for each medicine, its associated structured knowledge as the background context. We then apply the same synthesis pipeline as in FTK to produce hundreds of thousands of candidate items that cover clinical usage, contraindications, administration, and therapeutic indications. All candidates undergo machine screening through a combination of

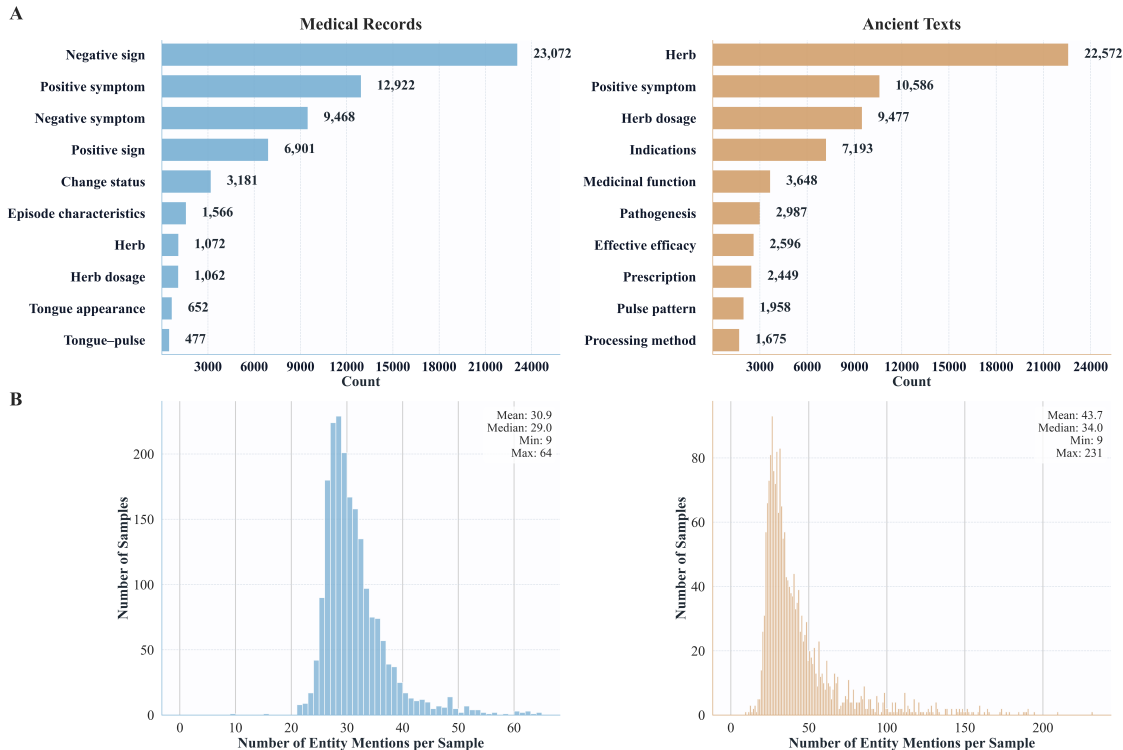


Figure 4: Statistical analysis of entity distribution in IE tasks. (A) Top 10 most frequently mentioned entity types in modern medical records (left, blue) and ancient medical texts (right, orange). Vertical bars show absolute mention counts per category. (B) Distribution of entity mentions per sample for medical records (left, blue) and ancient texts (right, orange).

regex-based heuristics (e.g., proportion of Latin characters, length constraints) and LLM-based quality scoring to remove redundant or low-quality items. Licensed TCM practitioners then review the remaining items to ensure linguistic clarity, terminological accuracy, and clinical plausibility. After multi-stage filtering and verification, 5,844 high-quality FTK items and 5,948 high-quality CPMK items are retained for inclusion in the benchmark.

3.3.4 Information Extraction Tasks

The information extraction task comprises two span-annotated corpora drawn from distinct text domains: de-identified modern Chinese electronic medical records (EMR) and classical TCM literature written in Literary Chinese. Following a human-machine workflow that integrated automatic NER suggestions with expert adjudication, we retained 2,000 samples for each corpus, yielding a total of 4,000 samples. The suite covers more than one hundred entity types, including positive and negative symptoms, physical signs, syndromes, Chinese herbs, formulas, condition changes, and episode characteristics. The clinical corpus contains 67 entity types with 11,886 mentions; the classical corpus contains 56 types with 19,876 mentions; in total there are 106 unique types. As shown in Figure 4 A, symptoms and physical signs are most frequent in EMR, reflecting the emphasis on symptom/sign documentation in clinical records; in the classical corpus, herbs and symptoms dominate, consistent with the prevalence of casebooks authored by renowned physicians. Figure 4 B presents the per-sample mention distribution: most samples contain 30–50 entities, and the densest classical case includes up to 231 entity mentions.

3.3.5 Diagnostic and Therapeutic Decision-Making Tasks

The auxiliary clinical decision component comprises two tasks: Diagnostic and Therapeutic Reasoning (DTR) and Decision Recognition (DR). From an initial pool of over 13,000 case records, we retained those with detailed documentation, specifically complete four-diagnoses notes and prescriptions containing at least five herbs, resulting in a curated DTR set of 2,000 high-quality cases. This set covers 1,562 distinct syndromes, 1,811 treatment principles, and 1,631 Chinese herbs. DTR integrates four interdependent subtasks that trace the clinical workflow: syndrome differentiation, treatment principle selection, prescription recommendation, and dosage estimation. Taken together, these subtasks provide a

comprehensive assessment of model competence in integrative clinical reasoning and actionable decision-making. Given that free-form prediction in DTR must search over extremely large label spaces conditioned only on the case text, we also construct DR by reframing each clinical problem as a single-choice question to measure cognitive alignment with expert reasoning while minimizing confounds from open-ended decoding. For DR, each case yields three items (syndrome, treatment principle, prescription), each with one correct option and four distractors, and is evaluated by accuracy. Option sets are built via embedding-based retrieval: we encode the query case with Qwen3-0.6B-Embedding, compute cosine similarity over a large indexed corpus, and collect the top 1,000 nearest neighbors. The correct option is the ground-truth label of the query case; distractors are sampled from labels observed among the neighbors with lower semantic similarity to the query, ensuring they remain contextually plausible yet discriminative.

In total, LingLan provides a comprehensive evaluation matrix integrating 44 indicators across 13 subtasks. This taxonomy ensures that every dimension of TCM knowledge and reasoning, across the spectrum from classical textual comprehension to modern clinical decision-making, is quantitatively represented.

3.3.6 Hard subset construction

To evaluate model performance under more challenging conditions, we construct a high-difficulty subset named LingLan-Hard. Item difficulty is estimated from per-instance scores across 14 models. For each item, we compute the across-model mean μ and variance σ , and define a composite difficulty score $D = (1 - \mu) + \lambda\sigma$, where $\lambda = 0.5$. Items are then ranked within each task in descending order of D , and the top 400 highest-difficulty items from each task are selected. In total, LingLan-Hard covers 13 tasks and contains 5,200 samples. The overall difficulty distribution across tasks is illustrated in *Appendix Figure 6*.

3.4 Evaluation Metrics

LingLan employs a unified metric matrix to ensure consistent, quantitative, and interpretable evaluation across heterogeneous task types. Each metric is aligned with the structural form of its corresponding task, allowing for fair comparison among classification, structured prediction, and regression subtasks.

3.4.1 Accuracy for Single- and Multiple-Choice

Accuracy is applied to discrete categorical tasks, including the TLE, FTK, CPMK, and DR subtasks. For single-choice questions, each instance has a unique correct option; for multiple-choice questions, an answer is considered correct only when all correct options are selected and no incorrect options are chosen:

$$\text{Accuracy} = \frac{N_{\text{correct}}}{N_{\text{total}}} \quad (1)$$

where N_{correct} and N_{total} denote the number of correctly answered and total questions, respectively.

3.4.2 Precision, Recall, and F1 for Cloze, Multiple-Choice IE, DTR, and DTR-F1

For any instance, let TP_i , FP_i , and FN_i denote true, false-positive, and false-negative counts as defined by the task-specific matching rule below. Precision, recall, and F1 are

$$\text{Precision}_i = \frac{TP_i}{TP_i + FP_i}, \quad \text{Recall}_i = \frac{TP_i}{TP_i + FN_i}, \quad \text{F1}_i = \frac{2 \text{Precision}_i \text{Recall}_i}{\text{Precision}_i + \text{Recall}_i} \quad (2)$$

and the dataset score is the macro average $\frac{1}{N} \sum_{i=1}^N \text{F1}_i$.

(1) Cloze (character-level matching). Given a predicted string \hat{s} and gold string s , treat them as multisets of characters. For each character c , let $n_{\hat{s}}(c)$ and $n_s(c)$ be their multiplicities. Define

$$TP_{\text{char}} = \sum_c \min(n_{\hat{s}}(c), n_s(c)), \quad FP_{\text{char}} = \sum_c \max(0, n_{\hat{s}}(c) - n_s(c)), \quad FN_{\text{char}} = \sum_c \max(0, n_s(c) - n_{\hat{s}}(c)) \quad (3)$$

Use the unified formulas above with $(TP_i, FP_i, FN_i) = (TP_{\text{char}}, FP_{\text{char}}, FN_{\text{char}})$. This captures partial correctness and orthographic variation at character granularity.

(2) Multiple-Choice. Each instance consists of selecting options from a predefined set. Let Y_i denote the gold set of correct options and \hat{Y}_i the predicted set. Define

$$TP_i = |\hat{Y}_i \cap Y_i|, \quad FP_i = |\hat{Y}_i \setminus Y_i|, \quad FN_i = |Y_i \setminus \hat{Y}_i|. \quad (4)$$

Apply the unified formulas. This regime measures exact set agreement between predicted and correct options.

(3) Information Extraction (instance-level multiset matching). Each instance is a multiset of typed surface forms over the universe U of pairs $u = (\text{type}, \text{text})$. Let G_i and \hat{G}_i be gold and predicted multisets with multiplicities $c_i(u)$ and $\hat{c}_i(u)$. Define

$$TP_i = \sum_{u \in U} \min(c_i(u), \hat{c}_i(u)), \quad FP_i = \sum_{u \in U} \max(0, \hat{c}_i(u) - c_i(u)), \quad FN_i = \sum_{u \in U} \max(0, c_i(u) - \hat{c}_i(u)) \quad (5)$$

where matching requires exact equality of entity type and normalized surface form. Apply the unified formulas.

(4) DTR. For syndrome differentiation, treatment principle, and prescription recommendation, each instance is a set of canonical labels without duplicates (order ignored). After schema-preserving normalization, counts are computed with a one-to-one matching rule that treats two labels as the same if either string contains the other. With gold set Y_i and predicted set \hat{Y}_i after schema-preserving normalization,

$$TP_i = |\hat{Y}_i \cap Y_i|, \quad FP_i = |\hat{Y}_i \setminus Y_i|, \quad FN_i = |Y_i \setminus \hat{Y}_i|. \quad (6)$$

Apply the unified formulas. This regime measures compositional correctness under label agreement.

(5) DTR-F1 (synonym-tolerant matching). To discount superficial lexical variation, align the predicted and gold sets by one-to-one bipartite matching at the instance level. For any predicted label \hat{y} and gold label y , compute the character-level F1 $F1_{\text{char}}(\hat{y}, y)$ as in the cloze definition. Construct a bipartite graph with edges where $F1_{\text{char}}(\hat{y}, y) \geq \tau$ (threshold $\tau = 0.7$) and take a maximum-cardinality matching M_i . Set

$$TP_i = |M_i|, \quad FP_i = |\hat{Y}_i| - |M_i|, \quad FN_i = |Y_i| - |M_i|, \quad (7)$$

then use the unified formulas. This preserves non-duplicated sets and enforces one-to-one alignment while tolerating near-synonymous surface forms.

3.4.3 Dosage Estimation

The dosage estimation task assesses the quantitative and proportional consistency between predicted and reference prescriptions. Given the vast combinatorial space of herbal prescriptions in TCM, achieving exact herb-level correspondence remains highly challenging even for advanced LLMs. To mitigate this difficulty, we introduce flexible alignment strategies by employing both inclusion-based matching (as used in DTR) and character-level F1 matching (as used in DTR-F1), thus ensuring that dosage evaluation remains fair and robust when minor lexical variations occur in herb names.

Each predicted herb is aligned with at most one reference herb according to the following criteria:

- (1) Inclusion-based matching: if the predicted herb name is contained within the reference herb name, or vice versa (e.g., “Ginseng” and “Ginseng Root”), the prediction is considered correct.
- (2) Character-level F1 matching: if inclusion does not hold, we compute the character-level F1 between the predicted and reference herb names as in the DTR-F1 evaluation; if it exceeds a threshold $\tau > 0.7$, the herb is also considered correct. This matching mechanism reduces the impact of superficial naming variation while maintaining one-to-one alignment between herbs, thus enabling consistent dosage evaluation.

Based on this aligned herb mapping, quantitative metrics are then computed to assess both absolute and proportional accuracy of dosage prediction.

(1) Mean Absolute Error (MAE). After alignment, the Mean Absolute Error quantifies the average deviation in dosage between matched herb pairs:

$$\text{MAE}_i = \begin{cases} \frac{1}{K_i} \sum_{k=1}^{K_i} |\hat{d}_{ik} - d_{ik}|, & K_i > 0, \\ 0, & \text{otherwise,} \end{cases} \quad (8)$$

where K_i denotes the number of matched herbs in instance i , and \hat{d}_{ik} and d_{ik} are predicted and reference dosages, respectively. MAE directly measures absolute quantitative accuracy. It provides a precise reflection of dosage deviation when most herbs are correctly matched between prediction and reference.

(2) Cosine Similarity. To better capture proportional dosage alignment, particularly when partial overlap exists between predicted and reference prescriptions, cosine similarity is also computed. Aligned dosage vectors $\hat{\mathbf{d}}_i$ and \mathbf{d}_i are constructed using inclusion- or character-F1-based matching; unmatched herbs are padded with zeros. For each instance:

$$\text{Cos}_i(\hat{\mathbf{d}}_i, \mathbf{d}_i) = \begin{cases} \frac{\hat{\mathbf{d}}_i^\top \mathbf{d}_i}{\|\hat{\mathbf{d}}_i\|_2 \|\mathbf{d}_i\|_2}, & \|\hat{\mathbf{d}}_i\|_2 > 0 \wedge \|\mathbf{d}_i\|_2 > 0, \\ 1, & \|\hat{\mathbf{d}}_i\|_2 = 0 \wedge \|\mathbf{d}_i\|_2 = 0, \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

Cosine similarity focuses on relative dosage proportions rather than absolute quantities. It remains informative even when some herbs are missing or mismatched, providing a complementary perspective on the model’s dosage reasoning ability.

Together, MAE and cosine similarity jointly characterize both absolute and proportional dimensions of dosage prediction, ensuring a comprehensive evaluation of quantitative reasoning in prescription generation.

3.4.4 Summary

LingLan establishes a unified and interpretable evaluation framework that ensures consistency across diverse task types in TCM. Accuracy is used for discrete-choice problems, instance-level or list-level Precision/Recall/F1 for structured and multi-label tasks, character-level F1 for short-text generation, and Cosine similarity together with MAE for quantitative dosage estimation. All metrics are computed in a macro-averaged manner and reported on both the full benchmark and its 400-item hard subsets, providing a comprehensive, fine-grained, and reproducible assessment of model performance across factual, reasoning, and quantitative dimensions.

3.5 Baseline Models and Evaluation Settings

To establish reliable baselines, we evaluated eleven representative LLMs spanning diverse architectures (dense vs. MoE), scales, and training paradigms (standard SFT/RLHF vs. explicit reasoning training). We restrict our evaluation to open-weight models; owing to privacy and data-governance constraints in medical scenarios, closed-source systems are not assessed in this release.

Qwen3 series. [45] Qwen3-4B, Qwen3-8B, Qwen3-14B, Qwen3-32B, Qwen3-30B-A3B, Qwen3-Next-80B-A3B-Thinking, and Qwen3-235B-A22B are members of the Qwen3 series, which includes dense and MoE variants with multilingual pretraining and an optional *thinking* mode for deliberative inference. We use official chat checkpoints where available and enable the family’s thinking capability for consistency with other reasoning-enabled baselines.

DeepSeek series. [7] DeepSeek-R1 is a reasoning-focused model trained via reinforcement learning to elicit multi-step latent reasoning without extensive supervised traces. DeepSeek-V3.1-Think is a reasoning-enabled variant based on the DeepSeek MoE architecture, evaluated in its think mode to standardize deliberative decoding across families.

Baichuan-M2-32B [46] is a 32-billion-parameter model developed on Qwen2.5-32B-Base. It incorporates a medical augmented-reasoning framework featuring a large-verifier mechanism and GRPO-style optimization to enhance stability and factual consistency in complex clinical reasoning tasks. The model integrates extensive medical knowledge during supervised fine-tuning and reinforcement alignment, aiming to improve interpretability and robustness in diagnostic reasoning.

GPT series [47] GPT-5 and GPT-5-mini are the latest models in the GPT series. GPT-OSS-20B and GPT-OSS-120B are open-weight reasoning models built on an efficient MoE architecture. They are trained via large-scale distillation and reinforcement learning, optimized for agentic capabilities such as research browsing, Python tool use, and function calling. Both models employ a structured chat format for robust instruction following and demonstrate strong performance across reasoning, coding, and safety benchmarks.

Chinese Medical and TCM LLMs [23, 22, 24, 21, 40] Since LingLan is a completely new dataset, and most of its data have not yet appeared in any internet corpus, it poses a high level of difficulty. At the current stage, LLMs focused on TCM have relatively small model sizes and suboptimal training methods, resulting in a significant gap compared with larger-scale models. Therefore, only a subset of models was evaluated, and the corresponding results are reported in the *Appendix Table 4*.

All models were evaluated under a zero-shot setting without any task-specific prompt engineering or fine-tuning. When supported, decoding hyperparameters were kept consistent: temperature $T = 0.6$, a maximum generation length of 8,192 tokens, and reasoning/thinking mode enabled during inference. For GPT models, we set

Table 2: LingLan results covering all subtasks. Each entry reports full-set / Hard-subset performance (left/right). Bold values denote the highest score for each subtask–metric. *Average* refers to the mean of all metrics, excluding MAE. Abbreviations: TLE (TCM Licensing Examination), FTK (Fundamental TCM Knowledge), CPMK (Chinese Patent Medicine Knowledge), IE (Information Extraction), DTR (Diagnostic & Therapeutic Reasoning), DR (Decision Recognition).

Dom.	Subtask	Metric	DeepSeek-R1	DeepSeek-V3.1	GPT-5	Qwen3-235B	Qwen3-32B	Baichuan-M2-32B	Qwen3-Next-80B	GPT-OSS-120B
TLE	Comprehensive	Accuracy	95.0 / 78.3	95.5 / 81.0	87.6 / 52.3	95.0 / 77.8	91.2 / 62.3	89.8 / 58.3	93.8 / 73.0	58.1 / 14.5
FTK	Single-choice	Accuracy	86.1 / 46.5	86.1 / 45.3	85.4 / 44.5	85.0 / 45.3	86.1 / 44.3	83.4 / 35.8	82.1 / 33.3	70.4 / 28.2
FTK	Multiple-choice	Accuracy	64.8 / 31.5	62.4 / 28.0	53.9 / 19.3	61.3 / 25.0	56.3 / 23.3	48.6 / 14.0	52.6 / 14.8	35.3 / 9.8
FTK	Multiple-choice	Precision	92.5 / 82.4	91.9 / 80.5	90.6 / 75.3	93.0 / 80.9	90.3 / 77.5	90.5 / 75.4	89.4 / 71.6	81.2 / 63.5
FTK	Multiple-choice	Recall	92.1 / 79.3	92.0 / 80.5	87.2 / 69.8	89.1 / 73.7	90.9 / 78.0	84.5 / 66.1	84.7 / 65.2	83.2 / 65.6
FTK	Multiple-choice	F1	91.0 / 77.8	90.6 / 77.5	87.3 / 69.3	89.5 / 73.7	89.1 / 74.8	85.3 / 66.4	84.7 / 64.2	80.0 / 61.3
FTK	Cloze	char-F1	58.2 / 20.9	58.8 / 19.7	54.9 / 19.8	59.9 / 23.5	56.1 / 17.5	50.2 / 16.2	56.0 / 15.9	42.1 / 14.0
CPM	Single-choice	Accuracy	74.3 / 31.0	73.3 / 25.5	67.3 / 18.0	68.6 / 17.8	66.9 / 16.3	60.5 / 12.8	63.1 / 8.0	44.5 / 13.0
CPM	Multiple-choice	Accuracy	48.0 / 24.5	42.2 / 15.0	35.4 / 6.5	44.8 / 15.3	35.5 / 6.5	33.2 / 5.5	34.2 / 3.8	24.7 / 4.0
CPM	Multiple-choice	Precision	86.4 / 73.5	84.3 / 69.6	82.4 / 62.7	86.1 / 69.7	81.9 / 62.1	83.1 / 63.5	80.6 / 58.7	76.0 / 53.9
CPM	Multiple-choice	Recall	90.6 / 76.3	91.5 / 78.4	87.8 / 72.2	87.7 / 71.1	87.8 / 70.0	82.2 / 59.9	84.1 / 62.0	82.5 / 61.8
CPM	Multiple-choice	F1	87.0 / 72.4	86.3 / 71.5	83.3 / 64.6	85.4 / 68.0	83.2 / 63.8	80.6 / 58.6	80.1 / 57.0	77.3 / 55.6
CPM	Cloze	char-F1	74.0 / 57.6	74.4 / 55.7	67.6 / 45.7	74.3 / 54.9	67.5 / 41.3	64.6 / 41.0	70.8 / 46.7	57.5 / 36.6
NER	Clinical EMR	Precision	66.6 / 46.7	64.6 / 46.7	72.8 / 56.2	69.1 / 50.5	73.3 / 61.3	61.8 / 45.5	38.1 / 17.3	59.7 / 41.4
NER	Clinical EMR	Recall	61.3 / 43.9	58.6 / 42.4	71.9 / 56.0	66.3 / 49.3	62.2 / 49.8	56.0 / 41.9	33.5 / 15.0	55.8 / 41.9
NER	Clinical EMR	F1	63.1 / 44.2	60.8 / 43.4	71.7 / 55.1	67.2 / 49.1	66.7 / 54.2	58.1 / 42.6	35.2 / 15.5	57.1 / 40.6
NER	Classical Texts	Precision	60.1 / 40.2	61.4 / 42.8	44.7 / 11.4	57.3 / 37.9	57.9 / 38.1	58.9 / 38.0	38.7 / 13.8	51.1 / 31.0
NER	Classical Texts	Recall	62.9 / 44.7	64.4 / 48.8	51.7 / 16.3	62.6 / 45.9	56.9 / 37.9	61.6 / 42.6	40.3 / 14.9	56.6 / 38.5
NER	Classical Texts	F1	60.8 / 41.5	62.1 / 44.7	47.5 / 13.1	59.1 / 40.4	56.6 / 36.9	59.5 / 39.2	39.0 / 13.8	53.0 / 33.5
DTR	Syndrome	Precision	9.9 / 1.3	11.0 / 1.5	5.0 / 0.7	9.0 / 1.2	7.7 / 0.6	9.1 / 1.2	10.7 / 1.5	3.2 / 0.6
DTR	Syndrome	Recall	13.3 / 1.7	14.3 / 1.8	8.4 / 1.1	13.4 / 1.9	11.2 / 0.9	12.6 / 1.6	12.6 / 1.7	4.5 / 0.7
DTR	Syndrome	F1	10.9 / 1.5	11.9 / 1.6	6.0 / 0.9	10.4 / 1.4	8.8 / 0.7	10.1 / 1.4	11.0 / 1.6	3.6 / 0.6
DTR	Treatment	Precision	14.2 / 1.9	14.3 / 0.8	10.5 / 2.1	12.9 / 0.9	11.5 / 1.3	11.7 / 0.6	13.6 / 0.6	9.6 / 0.9
DTR	Treatment	Recall	17.0 / 2.3	16.4 / 0.8	19.2 / 4.6	15.6 / 1.2	14.4 / 1.5	15.0 / 0.9	14.5 / 0.9	11.2 / 1.2
DTR	Treatment	F1	14.8 / 1.9	14.7 / 0.8	13.1 / 2.8	13.6 / 1.0	12.3 / 1.3	12.6 / 0.7	13.5 / 0.7	9.9 / 1.0
DTR	Prescription	Precision	35.7 / 19.7	36.9 / 19.9	27.6 / 17.1	33.4 / 19.0	31.0 / 17.0	32.1 / 18.0	32.4 / 17.4	30.2 / 17.3
DTR	Prescription	Recall	30.6 / 17.9	29.6 / 17.0	39.3 / 25.9	33.5 / 20.1	29.3 / 16.9	29.4 / 17.2	29.9 / 16.8	22.9 / 13.6
DTR	Prescription	F1	32.3 / 18.2	32.1 / 17.7	31.9 / 20.1	32.7 / 18.9	29.5 / 16.4	30.0 / 17.0	30.4 / 16.4	25.4 / 14.7
DTR	Dosage	MAE	4.1 / 4.0	4.2 / 4.1	4.5 / 4.9	4.2 / 4.2	4.4 / 4.7	4.2 / 4.1	4.0 / 3.9	4.4 / 3.7
DTR	Dosage	Cosine	30.9 / 16.2	31.0 / 16.4	31.4 / 19.7	31.2 / 17.7	29.0 / 15.5	28.7 / 15.4	29.1 / 15.2	23.3 / 12.8
DTR-F1	Syndrome	Precision	20.4 / 8.9	21.5 / 10.0	12.1 / 5.8	19.0 / 9.5	16.5 / 7.8	18.1 / 7.9	20.9 / 9.3	10.6 / 5.4
DTR-F1	Syndrome	Recall	26.6 / 11.2	27.1 / 11.2	20.1 / 9.1	27.4 / 12.3	23.8 / 10.5	24.8 / 10.3	24.1 / 9.7	14.9 / 6.8
DTR-F1	Syndrome	F1	22.0 / 9.6	22.9 / 10.1	14.6 / 6.9	21.6 / 10.4	18.7 / 8.7	20.0 / 8.5	21.3 / 9.1	11.9 / 5.8
DTR-F1	Treatment	Precision	21.0 / 8.1	21.7 / 6.2	15.0 / 5.7	19.8 / 6.4	18.0 / 6.7	18.0 / 5.1	21.1 / 5.9	14.6 / 4.3
DTR-F1	Treatment	Recall	25.2 / 10.6	24.9 / 7.8	28.1 / 12.3	24.1 / 8.7	23.0 / 9.0	23.3 / 7.3	22.5 / 7.2	17.7 / 5.9
DTR-F1	Treatment	F1	22.1 / 8.9	22.5 / 6.7	18.9 / 7.6	21.0 / 7.2	19.5 / 7.5	19.6 / 5.9	21.1 / 6.4	15.3 / 4.9
DTR-F1	Prescription	Precision	35.6 / 19.8	37.0 / 20.0	27.7 / 17.3	33.4 / 19.2	30.9 / 17.0	31.6 / 18.0	33.5 / 17.9	30.1 / 17.4
DTR-F1	Prescription	Recall	31.0 / 18.2	29.6 / 17.1	39.7 / 26.1	33.6 / 20.2	29.6 / 17.0	30.2 / 17.8	30.8 / 17.2	23.0 / 13.7
DTR-F1	Prescription	F1	32.4 / 18.4	32.2 / 17.8	32.0 / 20.3	32.7 / 19.0	29.6 / 16.5	30.2 / 17.3	31.2 / 16.9	25.5 / 14.8
DTR-F1	Dosage	MAE	4.1 / 4.0	4.2 / 4.2	4.5 / 5.0	4.2 / 4.2	4.4 / 4.7	4.2 / 4.2	4.0 / 3.9	4.4 / 3.8
DTR-F1	Dosage	Cosine	31.2 / 16.5	31.1 / 16.5	31.7 / 20.1	31.3 / 17.8	29.2 / 15.6	29.0 / 16.1	29.7 / 15.6	23.5 / 12.9
DR	Syndrome	Accuracy	86.7 / 39.8	85.6 / 35.5	86.9 / 41.5	85.4 / 40.8	86.0 / 39.3	83.6 / 30.5	84.9 / 32.8	77.7 / 34.8
DR	Treatment	Accuracy	78.7 / 22.8	77.6 / 20.3	78.1 / 22.0	79.0 / 20.8	80.0 / 22.3	76.8 / 15.3	77.6 / 17.8	75.3 / 31.8
DR	Prescription	Accuracy	89.1 / 51.2	90.0 / 55.3	90.6 / 56.8	88.7 / 49.3	87.5 / 45.3	87.9 / 44.0	88.4 / 47.0	78.5 / 37.3
Average			51.1 / 31.9	50.9 / 31.2	48.1 / 28.0	50.6 / 30.8	48.4 / 28.8	47.1 / 26.2	44.9 / 22.6	40.7 / 23.0

reasoning_effort="medium" to prevent frequent exceedance of the 8,192-token budget observed under the "high" setting, while maintaining sufficient reasoning depth.

4 Experimental Results

Under a unified zero-shot protocol, with identical decoding settings where supported, we evaluate 14 models on 13 subtasks spanning five domains. Results are reported for both the full test sets and the curated 400-item hard subsets, and all scores are macro-averaged within each subtask. The complete results are presented in Table 2 (models with parameters greater than or equal to 32B) and Table 3 (models with parameters below 32B). Overall averages on the full sets cluster in the high 50s (for example, DeepSeek-R1 51.1, DeepSeek-V3.1 50.9), while the hard subsets lower the averages to approximately 20, indicating a consistent degradation under increased difficulty. Figure 5 provides a more intuitive visualization of the performance differences across models on both the full and hard subsets. The overall average for each model is obtained by first averaging over sub-tasks and metrics (accuracy, F1, cosine, precision, recall) within each task type, then taking the mean of these per-task-type scores so that each task type is weighted equally.

Table 3: Small-scale LLM results covering all subtasks. Each entry reports full-set / Hard-subset performance (left/right). Bold values denote the highest score for each subtask–metric. *Average* refers to the mean of all metrics, excluding MAE. Abbreviations: TLE (TCM Licensing Examination), FTK (Fundamental TCM Knowledge), CPMK (Chinese Patent Medicine Knowledge), IE (Information Extraction), DTR (Diagnostic & Therapeutic Reasoning), DR (Decision Recognition).

Dom.	Subtask	Metric	Qwen3-32B	Qwen3-30B-A3B	Qwen3-14B	Qwen3-8B	Qwen3-4B	gpt-5-mini	gpt-oss-20b
TLE	Comprehensive	Accuracy	91.2 / 62.3	88.8 / 55.3	87.6 / 50.0	84.6 / 40.5	76.6 / 30.3	76.3 / 24.5	48.3 / 13.5
FTK	Single-choice	Accuracy	86.1 / 44.3	82.6 / 32.0	83.0 / 34.5	80.8 / 33.8	76.2 / 28.5	77.0 / 28.7	65.4 / 29.3
FTK	Multiple-choice	Accuracy	56.3 / 23.3	52.0 / 17.0	49.5 / 13.8	47.8 / 13.0	42.4 / 12.3	46.2 / 10.8	24.8 / 6.0
FTK	Multiple-choice	Precision	90.3 / 77.5	89.8 / 75.2	90.6 / 75.2	88.4 / 71.8	86.1 / 69.0	86.1 / 67.7	76.1 / 59.5
FTK	Multiple-choice	Recall	90.9 / 78.0	88.3 / 73.5	85.7 / 69.4	86.8 / 70.2	84.9 / 67.6	86.8 / 69.0	74.9 / 58.1
FTK	Multiple-choice	F1	89.1 / 74.8	87.4 / 71.0	86.3 / 68.7	85.8 / 67.6	83.5 / 64.8	84.8 / 65.3	72.9 / 55.0
FTK	Cloze	char-F1	56.1 / 17.5	53.9 / 17.9	54.1 / 14.3	51.1 / 12.3	49.0 / 15.9	49.2 / 17.9	35.9 / 10.6
CPM	Single-choice	Accuracy	66.9 / 16.3	58.4 / 5.5	58.9 / 9.8	56.6 / 9.0	51.3 / 9.3	51.8 / 9.5	37.6 / 11.8
CPM	Multiple-choice	Accuracy	35.5 / 6.5	34.7 / 5.8	34.4 / 5.3	31.4 / 7.2	28.7 / 4.5	28.9 / 4.0	18.5 / 3.8
CPM	Multiple-choice	Precision	81.9 / 62.1	81.8 / 60.6	82.5 / 61.5	81.3 / 61.6	79.3 / 59.3	78.2 / 54.9	72.6 / 51.1
CPM	Multiple-choice	Recall	87.8 / 70.0	85.4 / 65.1	84.9 / 64.6	84.2 / 66.3	84.7 / 67.7	82.7 / 57.9	75.4 / 55.4
CPM	Multiple-choice	F1	83.2 / 63.8	82.0 / 60.7	82.0 / 60.4	80.9 / 61.6	80.0 / 60.7	78.7 / 54.3	71.6 / 50.8
CPM	Cloze	char-F1	67.5 / 41.3	66.3 / 41.3	66.8 / 39.8	64.1 / 37.9	61.4 / 37.0	65.1 / 43.3	49.9 / 29.5
NER	Clinical EMR	Precision	73.3 / 61.3	69.3 / 57.2	74.6 / 62.8	67.7 / 55.1	64.7 / 52.7	52.1 / 35.2	55.3 / 39.5
NER	Clinical EMR	Recall	62.2 / 49.8	48.4 / 37.9	62.1 / 52.0	57.6 / 45.1	52.3 / 41.3	57.0 / 43.6	51.6 / 38.5
NER	Clinical EMR	F1	66.7 / 54.2	56.2 / 44.6	67.1 / 55.7	61.5 / 48.6	57.2 / 45.4	54.0 / 38.3	52.8 / 38.0
NER	Classical Texts	Precision	57.9 / 38.1	59.7 / 38.9	59.8 / 39.6	58.1 / 35.4	54.5 / 35.0	48.9 / 23.9	50.3 / 26.6
NER	Classical Texts	Recall	56.9 / 37.9	56.3 / 37.3	56.3 / 38.3	54.7 / 34.4	49.2 / 32.0	60.5 / 35.6	54.0 / 31.8
NER	Classical Texts	F1	56.6 / 36.9	57.2 / 37.0	57.3 / 38.0	55.6 / 34.0	51.0 / 32.4	53.4 / 28.0	51.3 / 28.1
DTR	Syndrome	Precision	7.7 / 0.6	8.4 / 1.4	7.6 / 1.2	7.5 / 1.0	6.9 / 1.4	3.4 / 0.3	3.3 / 1.4
DTR	Syndrome	Recall	11.2 / 0.9	12.9 / 2.0	11.3 / 1.2	9.9 / 1.1	9.8 / 1.7	5.6 / 0.4	5.4 / 2.1
DTR	Syndrome	F1	8.8 / 0.7	9.8 / 1.6	8.8 / 1.2	8.2 / 1.0	7.8 / 1.5	4.1 / 0.3	3.9 / 1.6
DTR	Treatment	Precision	11.5 / 1.3	13.0 / 0.8	11.0 / 0.8	12.0 / 0.8	11.6 / 1.3	13.4 / 1.3	6.5 / 0.8
DTR	Treatment	Recall	14.4 / 1.5	15.4 / 1.1	14.3 / 1.1	14.7 / 1.0	14.4 / 1.7	15.7 / 1.5	7.8 / 0.9
DTR	Treatment	F1	12.3 / 1.3	13.6 / 0.9	12.0 / 0.9	12.8 / 0.9	12.4 / 1.5	12.9 / 1.2	6.7 / 0.8
DTR	Prescription	Precision	31.0 / 17.0	32.4 / 18.1	31.1 / 17.2	31.7 / 17.1	30.4 / 16.7	27.6 / 16.1	24.8 / 15.7
DTR	Prescription	Recall	29.3 / 16.9	26.0 / 14.9	29.6 / 17.2	26.4 / 14.7	25.0 / 14.0	30.8 / 19.0	21.2 / 14.0
DTR	Prescription	F1	29.5 / 16.4	28.3 / 15.8	29.8 / 16.7	28.2 / 15.2	26.8 / 14.7	28.5 / 16.9	22.2 / 14.2
DTR	Dosage	MAE	4.4 / 4.7	4.4 / 4.0	4.2 / 4.2	4.3 / 4.1	4.1 / 4.1	5.3 / 5.3	4.5 / 4.3
DTR	Dosage	Cosine	29.0 / 15.5	28.1 / 15.3	28.6 / 15.7	27.2 / 14.2	26.3 / 14.5	27.9 / 16.0	20.0 / 12.4
DTR-F1	Syndrome	Precision	16.5 / 7.8	17.3 / 9.2	16.5 / 8.3	17.4 / 8.2	15.5 / 8.1	10.1 / 4.7	7.7 / 3.4
DTR-F1	Syndrome	Recall	23.8 / 10.5	26.5 / 13.0	24.1 / 11.3	22.5 / 10.2	21.9 / 10.7	16.0 / 6.4	11.9 / 4.7
DTR-F1	Syndrome	F1	18.7 / 8.7	20.1 / 10.4	18.8 / 9.2	18.8 / 8.8	17.5 / 8.8	11.9 / 5.3	8.8 / 3.8
DTR-F1	Treatment	Precision	18.0 / 6.7	20.1 / 7.7	16.8 / 6.2	18.8 / 6.1	17.4 / 6.2	17.2 / 4.0	10.3 / 3.8
DTR-F1	Treatment	Recall	23.0 / 9.0	23.9 / 10.2	22.4 / 8.8	23.4 / 8.5	21.8 / 8.6	22.6 / 7.1	13.4 / 5.4
DTR-F1	Treatment	F1	19.5 / 7.5	21.1 / 8.6	18.6 / 7.1	20.2 / 6.9	18.8 / 7.1	17.4 / 4.8	11.1 / 4.3
DTR-F1	Prescription	Precision	30.9 / 17.0	32.4 / 18.1	30.9 / 17.1	31.7 / 17.2	30.4 / 16.8	27.7 / 16.3	25.4 / 15.7
DTR-F1	Prescription	Recall	29.6 / 17.0	26.0 / 14.9	29.9 / 17.4	26.5 / 14.8	25.1 / 14.1	31.1 / 19.2	21.9 / 14.2
DTR-F1	Prescription	F1	29.6 / 16.5	28.3 / 15.8	29.8 / 16.8	28.3 / 15.3	26.9 / 14.8	28.7 / 17.1	22.9 / 14.4
DTR-F1	Dosage	MAE	4.4 / 4.7	4.4 / 4.0	4.2 / 4.3	4.3 / 4.1	4.1 / 4.1	5.3 / 5.4	4.6 / 4.4
DTR-F1	Dosage	Cosine	29.2 / 15.6	28.2 / 15.3	28.7 / 15.8	27.4 / 14.5	26.3 / 14.6	28.1 / 16.2	20.8 / 12.5
DR	Syndrome	Accuracy	86.0 / 39.3	84.1 / 29.5	84.2 / 34.5	82.8 / 29.0	77.3 / 25.3	80.7 / 26.0	70.3 / 30.3
DR	Treatment	Accuracy	80.0 / 22.3	80.4 / 23.8	78.4 / 19.0	78.6 / 20.3	74.8 / 24.0	78.1 / 26.0	68.2 / 33.5
DR	Prescription	Accuracy	87.5 / 45.3	86.7 / 40.5	87.4 / 43.3	84.9 / 35.3	82.4 / 36.0	84.4 / 39.5	64.5 / 25.5
Average			48.4 / 28.8	47.2 / 26.7	47.5 / 27.2	46.2 / 25.4	43.8 / 24.5	43.8 / 23.3	36.1 / 20.8

4.1 Overall Performance Across Domains

Licensing-style questions attain the highest accuracies on both the full sets and the hard subsets. In FTK and CPMK, single-choice accuracy is lower than per-option F1 on multi-choice, and cloze scores are lower than multi-choice; degradations on the hard subsets are larger for single-choice and cloze than for multi-choice. For information extraction, EMR scores exceed classical-text scores on both settings, and the gap persists on the hard subsets. In diagnostic–therapeutic reasoning, strict exact matching (DTR) yields lower scores than the synonym-tolerant variant (DTR-F1) across syndrome, treatment, and prescription, while reformulating the same problems as decision recognition

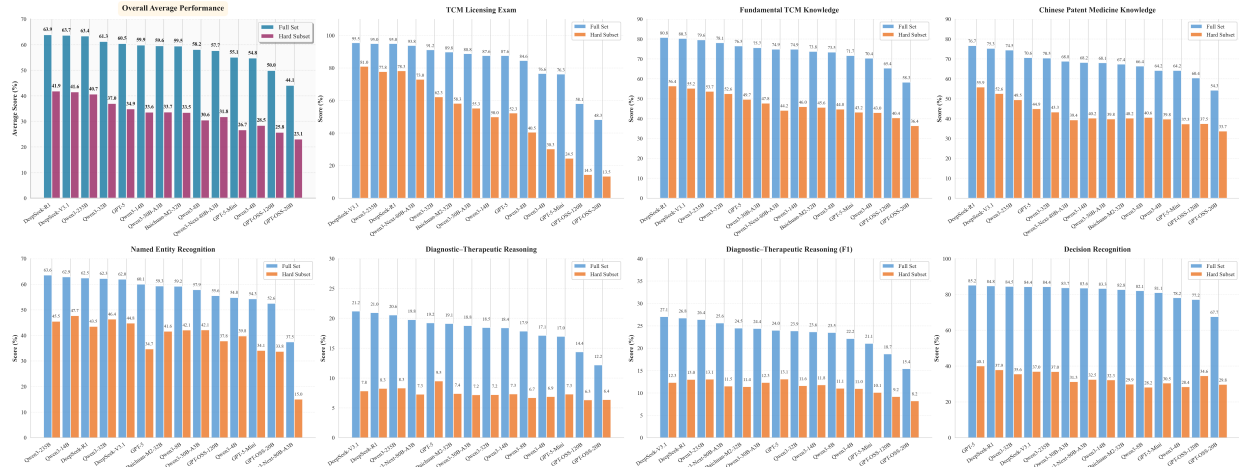


Figure 5: Model performance across all subtasks. Overall average refers to the mean performance across all tasks.

(DR) produces the highest accuracies among these formats. For dosage estimation, cosine similarity varies more across systems, whereas MAE remains within a comparatively narrow range on both the full and hard sets.

4.2 Overall Performance Across Models

On the full sets, overall averages cluster into three bands: a leading group centered near 51 (DeepSeek-R1, DeepSeek-V3.1, Qwen3-235B), a middle band in the high 40s (Qwen3-32B, GPT-5, Baichuan-M2-32B), and a lower band in the mid to low 40s (Qwen3-Next-80B, GPT-OSS-120B). On the hard subsets, the corresponding means concentrate near 31, with absolute declines averaging about 20 points across models. DeepSeek-V3.1 achieves the highest licensing-exam accuracy and the strongest CPM cloze on the full set; DeepSeek-R1 attains the highest multiple-choice F1 in FTK and CPM; Qwen3-235B leads FTK cloze; GPT-5 attains the top EMR extraction F1, the highest DR prescription accuracy, and the best dosage cosine.

4.3 Hard Subset Evaluation

All systems show pronounced significant on the hard subsets across domains, with the magnitude varying by task type and model. Among knowledge-oriented tasks, FTK multi-choice accuracy commonly declines by roughly 20–35 absolute points, and FTK cloze character-F1 drops by about 30–40 points. In CPMK, multi-choice F1 decreases by about 15–20 points, while CPMK cloze exhibits a smaller but still substantial contraction. Information extraction shows moderate-to-large declines: EMR F1 typically falls by 10–17 points, whereas classical-text F1 drops by about 15–30 points. In the diagnostic–therapeutic suite, strict DTR prescription F1 contracts from the low 30s to the high teens; dosage cosine similarity decreases by roughly 10–12 points. The decision-recognition reformulation (DR) remains the most accurate format on the full sets but also suffers the largest absolute drops in some subtasks, with syndrome accuracy falling from the high 80s to the low 40s, treatment often declining from around 80 to the low 20s, and prescription from about 90 to the mid-50s. Model-wise, high-capacity systems generally retain stronger hard-subset performance, yet all experience notable erosion. GPT-5 exhibits asymmetric robustness, leading EMR extraction but dropping sharply on classical-text extraction. Mid-tier and smaller models show steeper relative declines in several knowledge and extraction subtasks; for instance, GPT-OSS-120B and gpt-oss-20b register large losses on TLE and multi-choice QA. Taken together, the hard-subset results consistently stress models across knowledge recall, span extraction, multi-label clinical reasoning, and dosage proportionality, providing a stringent lens on robustness beyond full-set performance.

5 Discussion

5.1 Impact of Model Scale and Architecture

Table 3 summarizes results for small LLMs across 13 subtasks (below 32B parameters); these results enable task-wise analysis of how model scale and architecture affect performance. On knowledge-oriented tasks (TLE, FTK, CPMK),

performance increases clearly with scale: larger models better maintain stem–option alignment and terminology recall, especially in multi-choice settings, yielding higher and more stable option-level precision/recall/F1. This pattern suggests that “parameterized memory” and option calibration benefit from greater capacity.

In contrast, IE does not improve monotonically with size: smaller models sometimes exceed larger ones in macro-F1 on both EMR and classical-text corpora, indicating that extraction accuracy is sensitive to modeling behavior rather than raw capacity alone. Under the IE protocol (exact match on entity type, surface string with no synonym merging or normalization), larger models frequently exhibit benevolent normalization (standardizing terms, adjusting punctuation, consolidating near-duplicates) and confident inference (widening span boundaries or adding implicitly suggested entities). These behaviors reduce exact-match rates and can introduce type or polarity (positive/negative) mismatches, lowering precision and recall. Smaller models tend to copy spans verbatim and keep tighter boundaries, aligning more closely with the scoring criterion and yielding steadier macro-F1 despite their lower overall capacity.

Along the syndrome, treatment, prescription, dosage chain (DTR task), a common bottleneck persists across scales: enforcing global consistency over multi-label sets and exercising compositional control over prescriptions. Gains from scaling are modest; when reframed as single-choice DR, the search space contracts and size advantages re-emerge, yet all models still drop markedly on the hard subset, indicating robustness and transfer have not kept pace. Architecturally, dense models show better overall balance on knowledge recall and extraction, while MoE models tend to gain marginal advantages on discrete category decisions (e.g., syndrome, treatment method), likely from expert routing; these gains do not extend to strongly compositional or continuous-space tasks such as prescription composition and dosage ratios.

Overall, scale chiefly improves memory-driven, output-constrained formats, whereas architectural effects are task-dependent—benefiting categorical judgment but remaining insensitive to compositional structure. Continued progress on extraction and compositional reasoning will require tighter terminology and boundary standards, synonym-aware matching, and explicit mechanisms for chain consistency and prescription control.

5.2 Knowledge-Oriented Tasks

We summarize three observations characteristic of licensing-style (TLE) and foundational knowledge (FTK and CPMK) assessment in LingLan.

On the full sets, top systems approach ceiling accuracy on the TLE. A plausible explanation is that pretraining corpora may contain historical exam items or close paraphrases, which enhances recall of taught curricula and familiar test formats. Notably, this is the only LingLan component that could plausibly overlap with web-accessible material; all other datasets are sourced from non-web resources. Even so, the hard subset still depresses accuracy and exposes sensitivity to carefully engineered distractors and cross-chapter integrations. For multiple choice, a consistent gap appears between instance-level accuracy and per-option F1. Models often identify many correct options but miscalibrate the predicted subset size, leading to over-selection or under-selection. This highlights option calibration, rather than knowledge coverage alone, as a key challenge that LingLan makes explicit by reporting both metrics. Cloze questions are evaluated with character-level F1, which tolerates minor lexical and formatting variation relative to exact match, reducing spurious penalties from near-synonyms while better reflecting underlying knowledge. The metric remains deterministic and reproducible, supporting fair cross-model comparison.

Across these knowledge-oriented tasks, Chinese-centric model families (e.g., DeepSeek and Qwen) consistently occupy the top tier. This pattern is consistent with several language-domain factors: higher coverage of TCM terminology (materia medica, formula names, syndrome lexicon) in their pretraining data; tokenizers optimized for Chinese scripts that reduce sequence fragmentation; and instruction tuning that emphasizes Chinese exam formats and short-form factual recall. The advantage is most pronounced on TLE, FTK, and CPMK, which are dominated by recall and recognition, while it narrows in settings that demand span-faithful extraction or compositional generation.

5.3 Information Extraction Tasks

The information extraction task comprises two subtasks: NER on de-identified EMRs and on classical TCM texts in Literary Chinese.

5.3.1 Genre Effects and Scoring Considerations

Table 2 summarizes results on modern EMR and classical-text corpora and reveals clear genre effects alongside model-specific patterns. On EMR, precision typically exceeds recall, which is consistent with conservative span boundaries in semi-structured clinical prose. On classical texts, both precision and recall are generally lower than on EMR, reflecting the challenges posed by archaic vocabulary, elliptical syntax, and long-distance dependencies that complicate span delimitation and entity typing. The instance-level multiset scoring protocol, which matches typed surface strings while

preserving duplicates, exposes over-extraction and under-extraction that would otherwise be obscured by set-level deduplication.

5.3.2 GPT: Strong on EMR, Weak on Classical Texts

Across models, GPT-5 yields the best F1 on EMR, reflecting effective adaptation to contemporary clinical Chinese and stable span decisions under distribution shift. On classical texts, however, GPT-5 falls behind models such as DeepSeek-V3.1, which better handle archaic vocabulary and compact syntax. The Qwen3 family remains competitive across both genres and typically shows higher precision than recall on EMR, consistent with a cautious extraction style. Models with limited Chinese clinical and classical coverage, exemplified by GPT-OSS-120B, underperform on both corpora. Models whose pretraining is predominantly Chinese, such as DeepSeek and Qwen, tend to generalize better to classical Chinese, likely due to broader exposure to historical corpora, tokenization optimized for Chinese characters, and instruction tuning on Chinese sources.

5.3.3 Hard-Subset Effects and Methodological Directions

Degradation on the hard subsets is systematic and more pronounced for classical texts than for EMR, reflecting rarer terminology, compositional symptom phrases, and atypical surface forms that complicate span boundaries and type assignment. Improving robustness will likely require genre-aware tokenization and prompting for classical Chinese, span decision mechanisms that incorporate broader context without inflating boundaries, and post hoc consistency checks that reconcile entity inventories within long documents while balancing precision and recall.

5.4 Diagnostic and Therapeutic Reasoning Tasks

The diagnostic and therapeutic suite contains four interdependent subtasks in a free-form setting (DTR) and a recognition setting (DR) that reformulates the first three as single-choice questions.

5.4.1 Free-form Composition Remains Challenging

Across models, strict DTR scores are low in absolute terms and drop further on the hard subset, which indicates that open-set composition remains challenging. Under exact matching, prescription F1 peaks in the low thirties on the full set and falls to the high teens on the hard subset. Introducing the synonym-tolerant protocol yields modest but consistent gains, suggesting that a nontrivial share of errors arises from surface variation rather than conceptual mismatch. Error profiles show a stable precision–recall tradeoff: larger-capacity systems tend to retrieve more relevant labels, which improves recall but also introduces additional, partially plausible items that depress precision; smaller systems are more conservative, which stabilizes precision but limits coverage of multi-component prescriptions. Typical failure modes include under-selection of secondary herbs, confusion between closely related treatment principles, and incomplete alignment among syndrome, treatment, and prescription when long-range symptom relations are required.

Dosage estimation exhibits a different pattern. Cosine similarity, which emphasizes proportional agreement among dosed herbs, ranges in the low thirties on the full set and declines by roughly ten points on the hard subset. Mean absolute error remains around four to five grams per herb on both splits, with smaller values indicating better absolute calibration. The two metrics emphasize complementary behaviors. Cosine similarity is informative when overlap between predicted and reference herbs is limited, since it rewards correct dose ratios even with partial matches. MAE becomes more discriminative when most herbs are matched, since it measures absolute deviation. Observed errors concentrate on global scaling (for example, near-correct proportions but uniformly larger or smaller doses), unit normalization inconsistencies across sources, and occasional outlier herbs with atypical gram ranges that dominate the MAE.

5.4.2 Recognition Reformulation Improves Reliability

Recasting the same clinical problems as decision recognition substantially raises accuracy. With a constrained option set constructed from nearest-neighbor cases, models align more closely with expert labels for syndrome, treatment, and prescription. Hard-subset accuracy still drops, but the gap relative to DTR indicates that a significant portion of difficulty in free-form reasoning stems from the size of the hypothesis space and from calibration over multi-label outputs, rather than from a lack of stored knowledge alone. This pattern suggests practical routes to improvement, including candidate generation and re-ranking pipelines, structure-aware decoding that enforces compatibility among syndrome, treatment, and prescription, and training objectives that penalize over-selection and reward proportionally correct dosage vectors.

5.5 Terminology Normalization for Reliable TCM Evaluation

Existing studies show that general-purpose LLMs can align closely with expert judgments on TCM-specific tasks [17]; however, this evidence is narrow and depends on single-case, manual adjudication. A central impediment is terminology: pervasive synonymy, variant canonical forms, and heterogeneous granularity in syndrome labels, treatment principles, herb names, and processing states complicate both human and automated scoring. The field currently lacks a synonym and normalization resource broad enough to support reliable evaluation across these tasks. In LingLan, we partially mitigate lexical variability by using character-level F1 for surface-form tolerance and by reformulating open-ended outputs as constrained decision recognition, which improves comparability without overfitting to any single phrasing.

5.6 Future Work

LingLan is text-only and excludes core multimodal cues in TCM such as tongue images, facial and complexion signals, and pulse waveforms, which limits assessment of perception-to-reasoning competence. We plan a multimodal extension that pairs curated images and waveforms with expert annotations and schema-aligned labels, and introduces cross-modal retrieval, grounding, and decision tasks under unified metrics. This effort will include standardized acquisition protocols, privacy-preserving releases, and benchmark splits to enable reproducible evaluation of end-to-end multimodal TCM reasoning.

6 Conclusion

We introduce LingLan, a large-scale, expert-curated benchmark for TCM that spans five domains and thirteen subtasks under a unified evaluation protocol. The suite standardizes metrics across heterogeneous formats and provides hard subsets that probe fine-grained knowledge, distractor resistance, and compositional robustness. A zero-shot assessment of fourteen contemporary LLMs establishes the first broad, comparable baselines for TCM-oriented capabilities. Results indicate near-ceiling performance on licensing-style knowledge tasks, but information extraction on classical TCM texts remains a pronounced weakness, with substantially lower precision and recall than on modern EMRs. Despite persistent gaps in multi-step clinical reasoning, reframing open-ended decisions as recognition tasks substantially narrows these gaps, indicating that hypothesis-space size and calibration strongly influence performance. The results suggest that closing the gap will depend on models that integrate clinical structure with calibrated, domain-aware reasoning, moving beyond factual recall toward clinically faithful decision-making. By releasing standardized evaluation tools, we aim to catalyze reproducible research and more clinically faithful TCM modeling.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [5] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- [6] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- [7] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shitong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.
- [8] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan

- Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [9] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421, 2021.
- [10] Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao, Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, et al. Cmb: A comprehensive medical benchmark in chinese. *arXiv preprint arXiv:2308.08833*, 2023.
- [11] Junling Liu, Peilin Zhou, Yining Hua, Dading Chong, Zhongyu Tian, Andrew Liu, Helin Wang, Chenyu You, Zhenhua Guo, Lei Zhu, et al. Benchmarking large language models on cmexam-a comprehensive chinese medical exam dataset. *Advances in Neural Information Processing Systems*, 36:52430–52452, 2023.
- [12] Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 23(6):bbac409, 2022.
- [13] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. Toward expert-level medical question answering with large language models. *Nature Medicine*, 31(3):943–950, 2025.
- [14] Fenglin Liu, Hongjian Zhou, Boyang Gu, Xinyu Zou, Jinfa Huang, Jing Wu, Yiru Li, Sam S Chen, Yining Hua, Peilin Zhou, et al. Application of large language models in medicine. *Nature Reviews Bioengineering*, pages 1–20, 2025.
- [15] Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023.
- [16] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- [17] Yu Liu, Yishan Yuan, Keming Yan, Yuanyuan Li, Valeria Sacca, Sierra Hodges, Mattia Cannistra, Pauline Jeong, Jiani Wu, and Jian Kong. Evaluating the role of large language models in traditional chinese medicine diagnosis and treatment recommendations. *NPJ Digital Medicine*, 8(1):466, 2025.
- [18] Ethan Goh, Robert Gallo, Jason Hom, Eric Strong, Yingjie Weng, Hannah Kerman, Joséphine A Cool, Zahir Kanjee, Andrew S Parsons, Neera Ahuja, et al. Large language model influence on diagnostic reasoning: a randomized clinical trial. *JAMA network open*, 7(10):e2440969–e2440969, 2024.
- [19] Mickael Tordjman, Zelong Liu, Murat Yuce, Valentin Fauveau, Yunhao Mei, Jerome Hadjadj, Ian Bolger, Haidara Almansour, Carolyn Horst, Ashwin Singh Parihar, et al. Comparative benchmarking of the deepseek large language model on medical tasks and clinical reasoning. *Nature medicine*, pages 1–1, 2025.
- [20] Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023.
- [21] Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: Tuning llama model with chinese medical knowledge. *arXiv preprint arXiv:2304.06975*, 2023.
- [22] Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 19368–19376, 2024.
- [23] Yizheng Dai, Xin Shao, Jinlu Zhang, Yulong Chen, Qian Chen, Jie Liao, Fei Chi, Junhua Zhang, and Xiaohui Fan. Tcmchat: A generative large language model for traditional chinese medicine. *Pharmacological Research*, 210:107530, 2024.
- [24] Rui Hua, Xin Dong, Yu Wei, Zixin Shu, Pengcheng Yang, Yunhui Hu, Shuiping Zhou, He Sun, Kaijing Yan, Xijun Yan, et al. Lingdan: enhancing encoding of traditional chinese medicine knowledge for clinical reasoning tasks with large language models. *Journal of the American Medical Informatics Association*, 31(9):2019–2029, 2024.

- [25] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [26] Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. In *Advances in Neural Information Processing Systems*, 2023.
- [27] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*, 2019.
- [28] Suryani Lukman, Yulan He, and Siu-Cheung Hui. Computational methods for traditional chinese medicine: a survey. *Computer methods and programs in biomedicine*, 88(3):283–294, 2007.
- [29] Zhilin Song, Guanxing Chen, and Calvin Yu-Chian Chen. Ai empowering traditional chinese medicine? *Chemical science*, 15(41):16844–16886, 2024.
- [30] Wenjing Yue, Xiaoling Wang, Wei Zhu, Ming Guan, Huanran Zheng, Pengfei Wang, Changzhi Sun, and Xin Ma. Tcmbench: A comprehensive benchmark for evaluating large language models in traditional chinese medicine. *arXiv preprint arXiv:2406.01126*, 2024.
- [31] Ping Yu, Kaitao Song, Fengchen He, Ming Chen, and Jianfeng Lu. Tcmd: A traditional chinese medicine qa dataset for evaluating large language models. *arXiv preprint arXiv:2406.04941*, 2024.
- [32] Shufeng Kong, Xingru Yang, Yuanyuan Wei, Zijie Wang, Hao Tang, Jiuqi Qin, Shuting Lan, Yingheng Wang, Junwen Bai, Zhuangbin Chen, et al. Mtcmb: A multi-task benchmark framework for evaluating llms on knowledge, reasoning, and safety in traditional chinese medicine. *arXiv preprint arXiv:2506.01252*, 2025.
- [33] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- [34] Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.
- [35] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [36] Yan Cai, Linlin Wang, Ye Wang, Gerard de Melo, Ya Zhang, Yanfeng Wang, and Liang He. Medbench: A large-scale chinese benchmark for evaluating medical large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17709–17717, 2024.
- [37] Jack Gallifant and Danielle S Bitterman. Humanity’s next medical exam: Preparing to evaluate superhuman systems, 2025.
- [38] Suhana Bedi, Hejie Cui, Miguel Fuentes, Alyssa Unell, Michael Wornow, Juan M Banda, Nikesh Kotecha, Timothy Keyes, Yifan Mai, Mert Oez, et al. Medhelm: Holistic evaluation of large language models for medical tasks. *arXiv preprint arXiv:2505.23802*, 2025.
- [39] Zhe Wang, Meng Hao, Suyuan Peng, Yuyan Huang, Yiwei Lu, Keyu Yao, Xiaolin Yang, and Yan Zhu. Tcmeval-sdt: a benchmark dataset for syndrome differentiation thought of traditional chinese medicine. *Scientific Data*, 12(1):437, 2025.
- [40] Sibowei, Xueping Peng, Yifei Wang, Tao Shen, Jiasheng Si, Weiyu Zhang, Fa Zhu, Athanasios V Vasilakos, Wenpeng Lu, Xiaoming Wu, et al. Biancang: a traditional chinese medicine large language model. *IEEE Journal of Biomedical and Health Informatics*, 2025.
- [41] Yang Wu, Feilong Zhang, Kuo Yang, Shuangfang Fang, Dechao Bu, Hui Li, Liang Sun, Hairuo Hu, Kuo Gao, Wei Wang, et al. Symmap: an integrative database of traditional chinese medicine enhanced by symptom mapping. *Nucleic acids research*, 47(D1):D1110–D1117, 2019.
- [42] Lin Huang, Duoli Xie, Yiran Yu, Huanlong Liu, Yan Shi, Tielu Shi, and Chengping Wen. Tcmid 2.0: a comprehensive resource for tcm. *Nucleic acids research*, 46(D1):D1117–D1120, 2018.
- [43] Qunsheng Zou, Kuo Yang, Zixin Shu, Kai Chang, Qiguang Zheng, Yi Zheng, Kezhi Lu, Ning Xu, Haoyu Tian, Xiaomeng Li, et al. Phenonizer: A fine-grained phenotypic named entity recognizer for chinese clinical texts. *BioMed Research International*, 2022(1):3524090, 2022.
- [44] James E Harrison, Stefanie Weber, Robert Jakob, and Christopher G Chute. Icd-11: an international classification of diseases for the twenty-first century. *BMC medical informatics and decision making*, 21(Suppl 6):206, 2021.
- [45] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

- [46] Chengfeng Dou, Chong Liu, Fan Yang, Fei Li, Jiyuan Jia, Mingyang Chen, Qiang Ju, Shuai Wang, Shunya Dang, Tianpeng Li, et al. Baichuan-m2: Scaling medical capability with large verifier system. *arXiv preprint arXiv:2509.02208*, 2025.
- [47] Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025.

A Appendix

A.1 Evaluation on TCM LLM

We evaluate six models under a unified zero-shot protocol with identical decoding: a general-purpose baseline (Qwen3-14B); a Chinese medical model (Baichuan2-13B-Chat); two TCM-specific systems fine-tuned from Baichuan2-13B-Chat (TCMChat, Lingdan-13B-TCPM); and two TCM-specific systems fine-tuned from Qwen2.5 (BianCang-7B, BianCang-14B), with BianCang-7B supporting an 8,192-token maximum context length. Qwen3-14B attains the highest overall averages on both full and hard splits (Table 4) and leads NER on EMR and classical corpora. BianCang-7B and BianCang-14B are competitive on licensing-style assessment and decision recognition, especially on hard subsets, but trail Qwen3-14B on structured extraction and strict diagnostic–therapeutic composition. TCMChat and Lingdan-13B-TCPM show limited capability beyond single-choice QA, consistent with format-centric fine-tuning and the 4,096-token context limit inherited from Baichuan2-13B-Chat. Current TCM-specific models do not provide a dedicated reasoning mode, which further contributes to the substantial performance gap relative to the Qwen3 family. In dosage evaluation, near-zero MAE for some TCM models reflects failure to align herb names under inclusion-based matching and should not be interpreted as perfect accuracy; matching rules are detailed in the Evaluation Metrics section.

A.2 Difficulty Profiles and Hard-Subset Design

Figure 6 visualizes item difficulty by sorting each dataset’s instances by an empirical difficulty score (blue curve), marking the mean and median with dashed lines, and using the 400th item as the Hard cutoff. Licensing and single-choice knowledge tasks show steep early declines that indicate many easy recall items with a concentrated hard head, whereas multi-choice knowledge exhibits a more gradual slope and broader dispersion. Cloze tasks remain elevated longer before tapering, reflecting lexical sensitivity that keeps many items moderately difficult. In information extraction, classical texts sit consistently higher and decay more slowly than EMRs, revealing a heavier tail due to archaic vocabulary and elliptical syntax. In diagnostic–therapeutic reasoning, synonym-tolerant matching lowers the curve relative to strict scoring yet leaves a persistent hard tail driven by compositional prescription errors and dosage proportionality. Decision-recognition panels display stepped profiles, consistent with items that are either straightforward or distinctly challenging. Collectively, the Hard subsets lie above both mean and median across tasks, isolating the high-difficulty region while preserving each task’s characteristic sources of difficulty.

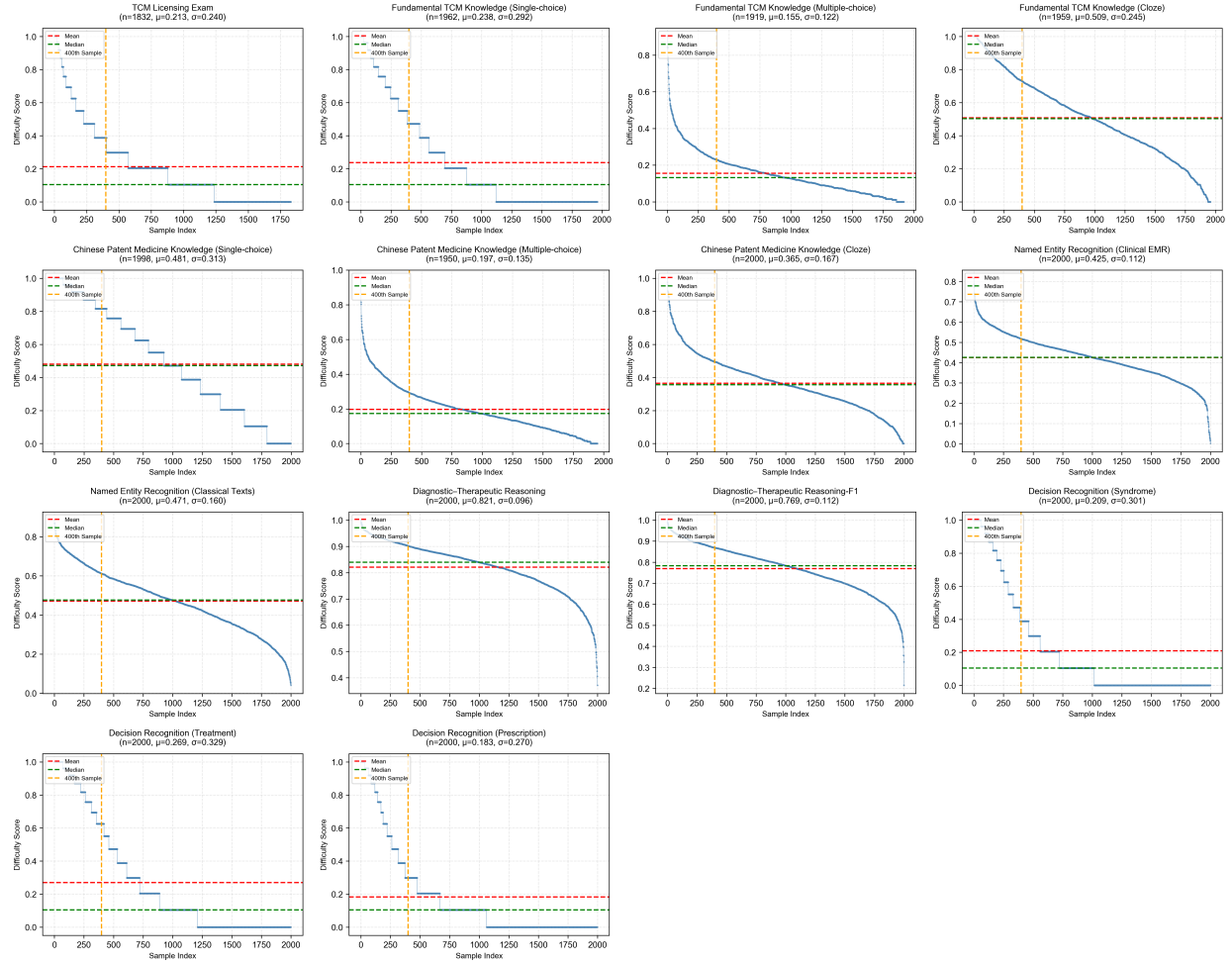


Figure 6: Difficulty distribution across all tasks.

Table 4: TCM LLM results covering all domains and subtasks. Each entry reports full-set / Hard-subset performance (left/right). Bold values denote the highest score for each subtask–metric.

Dom.	Subtask	Metric	Qwen3-14B	Baichuan2-13B-Chat	Lingdan-13B-TCPM	TCMChat	BianCang-7B	BianCang-14B
TLE	Comprehensive	Accuracy	87.6 / 50.0	50.1 / 21.8	44.3 / 20.8	31.6 / 19.5	88.9 / 70.0	91.2 / 73.5
FTK	Single-choice	Accuracy	83.0 / 34.5	60.6 / 28.5	48.4 / 26.3	11.1 / 7.5	78.2 / 43.3	82.6 / 49.3
FTK	Multiple-choice	Accuracy	49.5 / 13.8	23.4 / 12.8	14.3 / 9.5	0.6 / 0.5	36.9 / 19.5	23.7 / 13.0
FTK	Multiple-choice	Precision	90.6 / 75.2	71.8 / 62.0	53.6 / 48.8	58.9 / 48.2	78.9 / 70.1	65.2 / 59.8
FTK	Multiple-choice	Recall	85.7 / 69.4	84.4 / 79.8	56.5 / 56.4	36.3 / 29.1	78.4 / 72.7	53.1 / 48.7
FTK	Multiple-choice	F1	86.3 / 68.7	75.9 / 67.7	52.4 / 49.0	39.9 / 32.4	76.2 / 68.0	55.8 / 50.3
FTK	Cloze	char-F1	54.1 / 14.3	33.2 / 13.3	8.5 / 4.5	13.9 / 6.5	35.4 / 15.0	38.7 / 14.6
CPM	Single-choice	Accuracy	58.9 / 9.8	41.6 / 24.8	28.3 / 16.3	19.2 / 18.5	58.4 / 31.8	62.5 / 35.5
CPM	Multiple-choice	Accuracy	34.4 / 5.3	16.8 / 5.3	9.7 / 3.5	0.8 / 0.8	24.7 / 11.3	14.1 / 4.5
CPM	Multiple-choice	Precision	82.5 / 61.5	71.5 / 60.0	51.3 / 48.1	55.1 / 46.4	75.5 / 64.1	39.0 / 31.1
CPM	Multiple-choice	Recall	84.9 / 64.6	83.5 / 76.5	57.9 / 59.2	28.1 / 25.2	83.0 / 77.8	36.6 / 28.3
CPM	Multiple-choice	F1	82.0 / 60.4	75.2 / 65.4	52.3 / 51.2	33.8 / 29.7	77.1 / 68.3	36.3 / 27.8
CPM	Cloze	char-F1	66.8 / 39.8	37.2 / 26.4	11.7 / 9.3	13.8 / 10.6	54.4 / 37.5	58.4 / 42.2
NER	Clinical EMR	Precision	74.6 / 62.8	27.0 / 25.0	1.2 / 1.1	0.0 / 0.1	25.5 / 21.7	38.1 / 33.2
NER	Clinical EMR	Recall	62.1 / 52.0	25.0 / 22.6	0.8 / 0.9	0.0 / 0.0	15.0 / 11.5	26.1 / 21.5
NER	Clinical EMR	F1	67.1 / 55.7	25.1 / 22.7	0.8 / 0.8	0.0 / 0.0	16.8 / 12.9	29.5 / 24.8
NER	Classical Texts	Precision	59.8 / 39.6	32.3 / 21.5	5.6 / 2.8	1.5 / 1.8	46.1 / 25.2	55.3 / 39.0
NER	Classical Texts	Recall	56.3 / 38.3	27.8 / 18.1	2.9 / 1.5	0.2 / 0.2	44.9 / 26.9	46.0 / 32.4
NER	Classical Texts	F1	57.3 / 38.0	29.0 / 18.7	3.4 / 1.7	0.3 / 0.4	44.1 / 24.5	49.0 / 34.0
DTR	Syndrome	Precision	7.6 / 1.2	4.0 / 1.5	1.0 / 0.0	0.7 / 0.2	5.8 / 1.5	9.4 / 1.7
DTR	Syndrome	Recall	11.3 / 1.2	5.2 / 1.3	0.6 / 0.0	0.8 / 0.3	4.7 / 0.9	6.5 / 1.0
DTR	Syndrome	F1	8.8 / 1.2	4.3 / 1.3	0.7 / 0.0	0.7 / 0.3	4.9 / 1.1	7.3 / 1.2
DTR	Treatment	Precision	11.0 / 0.8	7.9 / 1.5	1.3 / 0.1	3.6 / 0.8	9.0 / 1.3	12.4 / 1.8
DTR	Treatment	Recall	14.3 / 1.1	7.6 / 1.4	0.9 / 0.1	2.6 / 0.5	6.3 / 1.0	7.8 / 1.0
DTR	Treatment	F1	12.0 / 0.9	7.4 / 1.3	1.0 / 0.1	2.8 / 0.6	7.1 / 1.1	9.1 / 1.2
DTR	Prescription	Precision	31.1 / 17.2	15.8 / 9.5	0.3 / 0.2	0.0 / 0.0	15.0 / 9.8	8.8 / 5.4
DTR	Prescription	Recall	29.6 / 17.2	12.7 / 8.1	0.2 / 0.2	0.0 / 0.0	15.4 / 10.1	8.3 / 5.6
DTR	Prescription	F1	29.8 / 16.7	13.6 / 8.3	0.2 / 0.2	0.0 / 0.0	14.4 / 9.1	8.2 / 5.1
DTR	Dosage	MAE	4.2 / 4.2	2.3 / 2.2	0.1 / 0.0	0.0 / 0.0	2.5 / 2.0	1.4 / 1.3
DTR	Dosage	Cosine	28.6 / 15.7	12.9 / 7.5	0.2 / 0.1	0.0 / 0.0	13.5 / 8.6	7.5 / 4.9
DTR-F1	Syndrome	Precision	16.5 / 8.3	8.8 / 5.5	1.6 / 0.1	0.8 / 0.5	12.8 / 7.8	19.9 / 8.4
DTR-F1	Syndrome	Recall	24.1 / 11.3	10.9 / 6.2	1.1 / 0.1	0.9 / 0.6	10.0 / 5.2	13.7 / 4.5
DTR-F1	Syndrome	F1	18.8 / 9.2	9.4 / 5.6	1.2 / 0.1	0.8 / 0.5	10.6 / 6.0	15.4 / 5.5
DTR-F1	Treatment	Precision	16.8 / 6.2	11.6 / 3.9	1.5 / 0.3	4.1 / 0.9	14.2 / 4.6	19.4 / 7.9
DTR-F1	Treatment	Recall	22.4 / 8.8	11.2 / 4.2	1.0 / 0.2	2.9 / 0.7	9.8 / 3.6	12.3 / 4.8
DTR-F1	Treatment	F1	18.6 / 7.1	10.9 / 3.9	1.1 / 0.2	3.1 / 0.8	11.0 / 3.8	14.3 / 5.7
DTR-F1	Prescription	Precision	30.9 / 17.1	20.1 / 13.0	1.9 / 0.9	2.0 / 0.6	17.5 / 10.8	23.6 / 14.8
DTR-F1	Prescription	Recall	29.9 / 17.4	14.5 / 9.5	1.5 / 0.7	1.1 / 0.3	18.3 / 11.7	23.2 / 15.4
DTR-F1	Prescription	F1	29.8 / 16.8	15.6 / 9.8	1.6 / 0.8	1.3 / 0.4	16.8 / 10.5	22.4 / 14.2
DTR-F1	Dosage	MAE	4.2 / 4.3	2.7 / 2.5	0.4 / 0.2	0.0 / 0.0	3.0 / 2.5	3.8 / 3.6
DTR-F1	Dosage	Cosine	28.7 / 15.8	14.8 / 9.0	1.4 / 0.6	0.1 / 0.0	15.9 / 9.9	20.9 / 13.0
DR	Syndrome	Accuracy	84.2 / 34.5	64.6 / 32.0	32.5 / 17.3	47.2 / 28.5	73.8 / 34.3	80.7 / 41.5
DR	Treatment	Accuracy	78.4 / 19.0	62.6 / 28.5	38.5 / 20.5	48.3 / 27.5	72.0 / 32.3	72.7 / 26.5
DR	Prescription	Accuracy	87.4 / 43.3	36.2 / 25.3	18.1 / 15.8	23.0 / 18.5	65.6 / 39.8	80.5 / 48.0
Average			47.5 / 27.2	30.3 / 20.7	14.6 / 11.2	11.7 / 8.6	35.3 / 23.7	33.5 / 21.5