

CCL25-Eval任务9系统报告：一种面向中医辨证与处方生成任务的检索增强大模型方法

康益扬, 姚佳琪, 吕腾啸, 徐博, 罗凌, 孙媛媛, 林鸿飞
大连理工大学计算机科学与技术学院
{kiang0920,1741917591,tengxiaolv}@mail.dlut.edu.cn
{xubo,lingluo,syuan,hflin}@dlut.edu.cn

摘要

本文面向CCL25-Eval任务9中的中医辨证辨病与中药处方推荐两个子任务，提出了一套基于大语言模型的系统性方法。在子任务1中，本文基于QLoRA方法对Qwen2.5-7B、Mistral-7B和Baichuan2-7B三种预训练模型进行高效微调，并引入多模型集成投票策略。在子任务2中，本文设计了融合向量检索、监督微调与强化学习的中药推荐框架，通过相似度检索构建候选处方集合，并利用强化学习优化模型的生成能力。最终在评测中获得总分0.5171（Task1得分0.5710，Task2得分0.4632），排名第四，验证了所提方法的有效性与实用性。

关键词： 中医辨证论治；大语言模型；参数高效微调；检索增强；强化学习

SystemReportforCCL25-EvalTask9:A Retrieval-Augmented Large Language Model Approach for TCM Syndrome Differentiation and Prescription Generation.

Yiyang Kang, Jiaqi Yao, Tengxiao Lü, Bo Xu, Ling Luo, Yuanyuan Sun, Hongfei Lin
School of Computer Science and Technology, Dalian University of Technology, Dalian
{kiang0920,1741917591,tengxiaolv}@mail.dlut.edu.cn
{xubo,lingluo,syuan,hflin}@dlut.edu.cn

Abstract

This paper proposes a systematic method based on a large language model for the two subtasks of TCM syndrome differentiation and TCM prescription recommendation in CCL25-Eval Task 9. In subtask 1, this paper efficiently fine-tunes the three pre-trained models of Qwen2.5-7B, Mistral-7B and Baichuan2-7B based on the QLoRA method, and introduces a multi-model ensemble voting strategy. In subtask 2, this paper designs a TCM recommendation framework that integrates vector retrieval, supervised fine-tuning and reinforcement learning, constructs a candidate prescription set through similarity retrieval, and uses reinforcement learning to optimize the model's generation ability. Finally, the total score of 0.5171 (Task1 score 0.5710, Task2 score 0.4632) was obtained in the official evaluation, ranking fourth, verifying the effectiveness and practicality of the proposed method.

Keywords: TCM Syndrome Differentiation , Large Language Model , Parameter-Efficient Fine-tuning , Retrieval-Augmented Generation , Reinforcement Learning

©2025 中国计算语言学大会
根据《Creative Commons Attribution 4.0 International License》许可出版

1 引言

中医辨证论治是中华传统医学的核心理念之一，强调“因人、因时、因地制宜”的个体化治疗原则。通过望、闻、问、切四诊收集临床信息，医生需对病人的主诉、体征、舌苔、脉象等多维症状进行归纳分析，进而辨识证型与疾病，制定对应的治疗方案与中药处方。这一过程不仅依赖于丰富的医学知识，还高度依赖临床经验与综合判断能力，因而在自动化建模与智能诊疗系统构建中具有很高的挑战性与研究价值。

为推动中医领域的智能诊疗研究，本次评测任务构建了包含1500条脱敏病历数据的数据集，覆盖10种证型、4种疾病和381种中药，设置“中医多标签辨证辨病”与“中药处方推荐”两个子任务，子任务1定义为给定一段患者的详细情况描述(包括病史、主诉、四诊信息等)，来预测出患者所患的疾病和证型；子任务2定义为给定患者的详细情况描述(包括病史、主诉、四诊信息等)，来推荐当前患者的中药处方。本文针对该任务展开了系统研究，在子任务1中，本文采用QLoRA 技术对Qwen2.5-7B、Mistral-7B 和Baichuan2-7B 等模型进行微调，并通过多模型集成投票提升辨病辨证的准确性；在子任务2中，本文提出了一种结合检索、监督微调和强化学习的中药处方推荐方法，先通过相似度检索获取中药推荐列表，接下来通过对大模型进行监督微调和强化学习，并多次采样，最终将得到的结果与检索阶段的结果整合重排。本文的贡献在于：

1. 提出了一种针对中医多标签辨证辨病任务的集成式建模框架，有效提升了模型的预测性能。
2. 设计了一种结合检索增强和强化学习的中药处方推荐方法，能够生成更具可靠性和准确性的中药推荐结果。

2 相关工作介绍

随着人工智能的快速发展，许多基于语料进行预训练的大语言模型也随之出现，例如OpenAI推出的ChatGPT(Kasneci et al., 2023)。大语言模型 (Large Language Model, LLM) 是指通过在海量语料数据上训练得到的深度神经网络模型，具备强大的语言建模和语义理解能力，通常包含数百亿个参数。大语言模型通常采用基于Transformer的架构(Vaswani et al., 2017)，通过自回归或自编码的方式预测文本序列中词语的分布，进而实现语言理解与生成任务。大语言模型在自然语言理解(Hendrycks et al., 2020)、文本生成(Touvron et al., 2023)和推理(Zhang et al., 2022)等方面表现出强大的能力，也逐渐成为构建复杂语言任务解决方案的核心技术。在医学领域，大语言模型也同样展现出了良好的应用前景，包括辅助诊断治疗(Yang et al., 2024b)、医学成像(Liu et al., 2023)、健康管理(Jin et al., 2024)和药物开发(Zheng et al., 2024)等场景。

2.1 中国传统医学及智能应用

中国传统医学 (Traditional Chinese Medicine, TCM) 是中华民族几千年医疗实践与哲学思想积淀而成的医学体系，具有独特的理论结构和诊疗方法。其核心强调“天人合一”、“阴阳五行”与“辨证论治”，通过望闻问切合参，全面获取病人的生理与病理信息，进而归纳为特定的“证型”，制定个体化的治疗方案(Unschuld, 2010)。与现代西医强调病因定位和标准化治疗不同，中医更重视个体差异与整体调和，注重疾病与情志、体质、环境的相互关系，其诊断与处方往往依赖经验积累和复杂的知识体系。在临床实践中，“辨证论治”是中医的核心诊疗原则，即通过综合分析症状、舌象、脉象等信息，辨识病因、病位、病性等因素，并据此归纳出证型，再选择合适的中药方剂予以治疗(Zhou et al., 2016)。例如，同一种疾病在不同体质患者中可能表现出不同的证型，其治疗方案亦可能截然不同，这种灵活性使中医更具个体化医疗的潜力。

近年来，随着中医现代化进程推进，越来越多研究致力于通过人工智能等手段重构中医诊疗过程，提升其规范性、智能化水平并增强其在全球健康体系中的应用价值。例如，Zhongjing(Yang et al., 2024b)是一个基于LLaMA通过专家反馈和多轮真实对话进行强化学习的中医大模型。Qibo(Zhang et al., 2024)是一个结合连续的预训练和监督微调，采用两阶段训练方法得到的中医大语言模型，旨在解决中医理论与现代医学之间的差异以及专业语料库匮乏等挑战。这些研究表明，中医与人工智能的融合正在逐步深化，为中医知识的结构化表达与

智能化应用奠定基础。随着中医语料资源的扩展与大模型能力的提升，人工智能算法将在中医诊断与治疗辅助中展现更广阔的应用前景。

2.2 医学药物推荐

药物推荐系统在现代医疗体系中扮演着关键角色，旨在基于患者的临床信息提供个性化、精准的用药方案，辅助医生进行高效、安全的临床决策，来降低可能潜在存在的用药风险(Su et al., 2020)。随着电子健康记录(Electronic Health Records, EHR)在医疗机构中的广泛部署，海量结构化与非结构化数据为个性化药物推荐提供了丰富的数据基础(Henry et al., 2016)。在多疾病共存、用药风险较高的情境中，推荐系统能够帮助医生发现潜在的药物相互作用、优化药物组合，有效提升治疗效果并降低临床错误率。然而，如何全面整合和建模患者的体征、诊断信息、实验检查和病程记录，仍是药物推荐系统面临的重要挑战(Ali et al., 2023)。

现有研究中，主流方法涵盖协同过滤、图神经网络(GNN)、序列建模与强化学习等技术路径。早期方法多依赖规则驱动或静态编码策略，仅基于诊断码和手术记录进行推荐，难以捕捉病情演变的动态特征。随着深度学习的发展，研究者逐步引入循环神经网络(RNN)、注意力机制和图卷积网络来建模患者病史与药物之间的复杂依赖关系。例如，SafeDrug模型利用药物分子结构构建双图编码器，以提高推荐的用药安全性(Yang et al., 2021)；GAMENet通过联合建模健康记录与药物知识图谱，有效融合纵向病史与药物交互知识(Shang et al., 2019)。然而，上述方法通常依赖复杂的特征工程与编码流程，限制了其在真实临床场景中的可扩展性与应用普适性。近年来，大语言模型的崛起为药物推荐任务提供了新思路。已有研究初步探索将大语言模型与蒸馏技术结合，用于迁移知识至轻量模型以进行推荐(Liu et al., 2024)，表明大语言模型具备对复杂医疗文本理解与推理的潜力，但其在药物推荐中的应用仍处于起步阶段，亟需进一步挖掘其在多模态医疗数据上的泛化能力与可解释性。

3 模型与方法

3.1 子任务1

中医多标签辨证辨病任务要求基于给定的患者临床文档，判断患者所患的证型和疾病。该任务旨在借助自然语言处理技术从病历和症状描述中提取信息，快速分析主证、兼证和疾病的关系，提高辨病辨证的效率和准确性，辅助医生更精准地完成诊断。子任务1输入包含：“ID”、“性别”、“职业”、“年龄”、“婚姻”、“病史陈述者”、“发病节气”、“主诉”、“症状”、“中医望闻切诊”、“病史”、“体格检查”、“辅助检查”；输出包括：“证型”、“疾病”，输出为预测患者所患的“疾病”和“证型”，如图1所示：

Input
"ID": "35"
"性别": "女"
"职业": "退休"
"年龄": "66岁"
"婚姻": "已婚"
"病史陈述者": "本人"
"发病节气": "立夏"
"主诉": "发作性胸闷20年，加重伴胸痛3月余"
"症状": "胸部疼痛，呈针刺样，胸闷不舒，心慌不安，气短乏力，眼干眼涩，口干口苦，纳可，食后反酸烧心，眠可，二便调。"
"中医望闻切诊": "中医望闻切诊：表情自然，面色暗红，形体正常，动静姿态，语气低，气息平；无异常气味，舌暗红、苔黄腻，舌下络脉曲张，脉弦。"
"病史": "现病史：患者20年前因劳累后出现胸闷，无胸痛，于***就诊，行心电图、心脏彩超等检查，诊断为“冠心病”，具体治疗不详，好转后出院。院外规律服用“拜阿司匹林”等，症状控制可。3月前劳累后出现上述症状加重，自服“复方丹参滴丸”……"
"体格检查": "体温：36.5℃ 脉搏：61次/分 呼吸：18次/分 血压：158/71mmHg (R)、152/71mmHg (L) Padua评分：3分(低危) 生命体征一般情况：患者老年，女，发育正常，营养良好，神志清楚，步入病房，查体合作……"
"辅助检查": "2020-4-29 冠脉CT示：LM轻度狭窄，LAD中度狭窄，LCX轻度狭窄，RCA轻度狭窄，PDA中度狭窄。(于***) 2020-5-12 心电图示：窦性心律，ST-T改变。"
Output
任务1:
"疾病": "胸痹心痛病"
"证型": "气虚血瘀证热蕴蕴证"
任务2:
"处方": "[丁香, 广藿香, 黄芪, 檀香, 砂仁, 木香, 草豆蔻, 附片 ……]"

图 1: 任务输入输出示例

针对子任务1，本文基于预训练大语言模型设计了一套集成式建模框架，旨在提升模型对中医临床文本中复杂症状信息的理解能力与标签预测准确性。在模型选型方面，本文选取了

三种主流的开源大语言模型作为基础模型，分别为Qwen2.5-7B(Yang et al., 2024a)、Mistral-7B(Jiang et al., 2023)和Baichuan2-7B(Yang et al., 2023)，并采用QLoRA (Quantized Low-Rank Adaptation) 方法对其进行轻量化微调(Hu et al., 2022)。QLoRA 是一种参数高效的微调方法，结合了4-bit 权重量化与低秩权重插入技术，能够在资源受限的环境下实现对大模型的高效训练，并保持与全参数微调接近的性能表现。

3.1.1 训练阶段

在训练阶段，本文采用Prompt-Response格式组织训练样本，其中输入为拼接后的患者临床信息（包括主诉、现病史、舌脉象等），输出为对应的证型标签和疾病标签。所有模型分别在任务所给数据量为800的训练集上独立微调，以增强模型间的多样性，为后续集成提供基础。在独立微调的实验过程中所使用的超参数如表1所示：

Learning Rate	Optimizer	LoRA Rank	LoRA Alpha	LoRA Dropout
2×10^{-5}	AdamW	64	16	0.05

表 1: 微调超参数设置

3.1.2 推理阶段

在推理阶段，采用多模型集成投票策略对三个微调模型的预测结果进行融合。具体而言，对于每条测试样本，分别获取三个模型独立预测的证型与疾病标签集合，并对每个标签统计投票次数，最终选择出现次数最多的标签作为模型的最终输出。该策略能够充分利用各模型在不同样本上的优势，提升整体预测的鲁棒性与稳定性。推理流程图如图2所示：

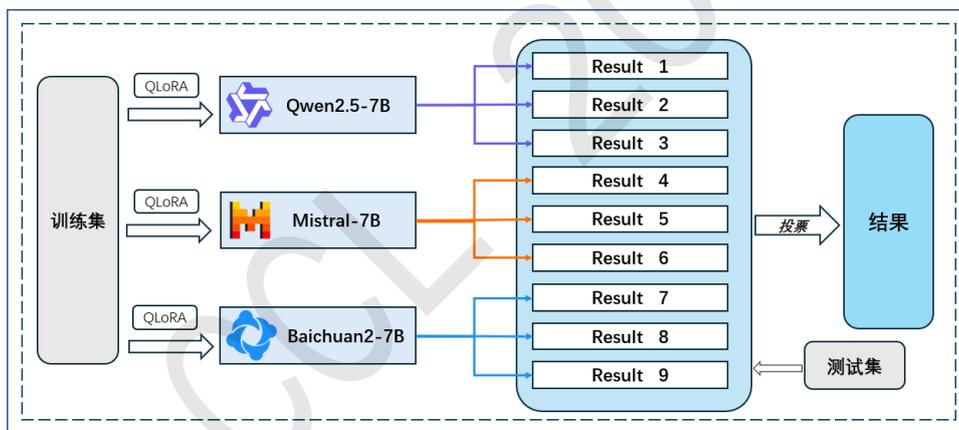


图 2: 推理集成框架图

3.2 子任务2

中药处方推荐任务要求基于给定的患者临床文档，为患者推荐合适的中药处方。该任务旨在根据患者的详细情况描述，自动推荐一组中药作为处方，如图1所示。在子任务2中，本文提出了一种结合检索增强、监督微调和强化学习的中药处方推荐方法。首先，通过检索模块，利用bge模型对测试集输入进行嵌入，并基于余弦相似度从训练集中检索出60条最相似的数据，提取其对应的中药处方，汇总形成中药列表。接着，在监督微调和强化学习模块中，以Qwen2.5-7B模型为基础，采用LoRA技术进行微调以提升模型性能，并进一步应用GRPO强化学习算法，结合与任务评估指标一致的奖励函数，优化模型的中药推荐策略。最后，在结果采样与重排阶段，使用强化学习优化后的模型对测试集输入进行60次采样，生成中药列表，并与检索模块的中药列表合并，通过频率统计选取出现频率最高的前16种中草药作为最终的中药处方输出。图3是该任务的整体框架图。

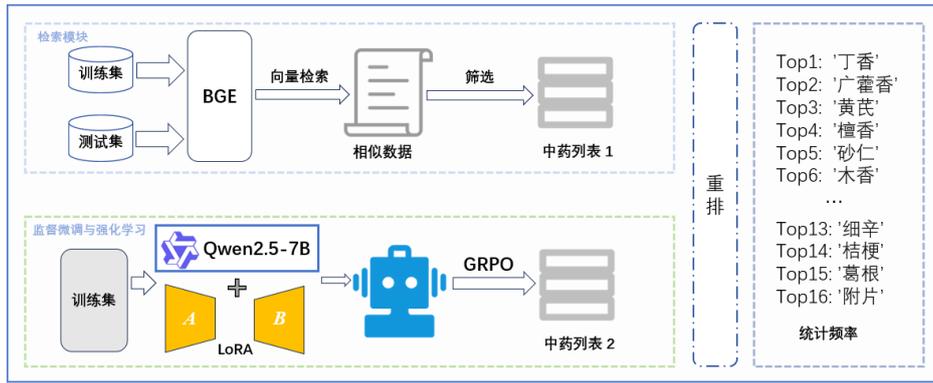


图 3: 子任务2整体框架图

3.2.1 检索模块

检索模块的目的是从训练集中找到与测试集输入相似的样本，以获取初步的中药推荐列表。具体步骤如下：

1. 相似性检索：对于测试集中的每一条输入，采用bge模型进行嵌入，并使用余弦相似度进行检索，在训练集中检索出与之最相似的60条数据。
2. 中药列表提取：从检索到的60条相似数据中提取对应的中药处方，将这些中药处方中的中草药进行汇总，形成检索阶段的中药列表。

3.2.2 监督微调与强化学习模块

为了进一步优化中药推荐结果，本阶段采用模型监督微调和强化学习相结合的方法。

1. 监督微调：选择Qwen 2.5 7B 模型作为基础模型，并使用训练集数据对其进行LoRA（Low-Rank Adaptation）微调。LoRA 是指通过在模型的权重矩阵中添加低秩矩阵来实现对模型的适应性调整，能够在不大幅增加计算成本的情况下提升模型的性能。
2. 强化学习：对微调后的模型应用GRPO（Generalized Reward Policy Optimization）强化学习算法。在强化学习过程中，设计了与子任务2评估指标相同的奖励函数，以指导模型学习更优的中药推荐策略。通过强化学习，模型能够根据奖励信号不断调整其输出，从而提高推荐结果的质量。

3.2.3 结果采样与重排

利用经过强化学习优化后的模型对测试集输入进行采样，采样次数为60次，每次采样生成一个中药列表。将采样得到的中药列表与检索模块提取的中药列表合并，形成一个综合的中药列表。对综合中药列表中的中草药进行频率统计，选取出现频率最高的前16种中草药作为最终的中药处方输出。

4 实验结果与分析

表 2 展示了本文的方法在验证集和测试集上的具体得分情况：

数据集	Task1-Score	Task2-Score	总分 (Overall Score)	排名
Baseline	0.4500	0.4100	0.4300	第68名
验证集	0.5650	0.4705	0.5177	第5名
测试集	0.5710	0.4632	0.5171	第4名

表 2: 在验证集和测试集上的评测结果

从表 2 可以看出, 本文提出的多模型集成方法在任务1 (辨证辨病) 中, 在验证集与测试集上均显著优于Baseline, Task1-Score 提升超过0.11。在任务2 (中药推荐) 中, 融合检索增强与强化学习策略的方案也取得了优于Baseline 的结果, 展现出良好的泛化能力。

为了进一步验证各模型对性能的影响, 本文在任务1上开展了多种单模型实验, 结果见表 3。其中, Mistral-7B 在微调后取得了最优的单模型结果 (0.5475), Qwen2.5-7B 与Baichuan2-7B 紧随其后, 三者性能表现相对接近。DeepSeek-7B 的得分略低, 但仍超过0.5, 说明在生成式任务中, 高质量的开源大模型普遍具备良好的中医文本建模能力。

相较之下, 将任务1转化为多标签分类任务并使用Roberta 进行微调的方案表现明显较差, 反映出模型参数规模和语言建模能力在中医诊疗类任务中的重要性。大型生成式语言模型在理解复杂症状描述和处理中文医学术语方面展现出更强的能力, 进一步验证了大模型在中医智能诊断中的应用潜力。

方法	Task1-Score
Qwen2.5-7B (生成式模型)	0.5450
Roberta (分类模型)	0.4650
Mistral-7B (生成式模型)	0.5475
DeepSeek-7B (生成式模型)	0.5050
Baichuan2-7B (生成式模型)	0.5425
Qwen2.5-7B (任务1+2联合生成)	0.5325
多模型集成	0.5650

表 3: 任务1不同模型在验证集上的表现

在任务2中 (表 4), 若直接基于Qwen2.5-7B 进行中药处方生成而不引入检索机制, 其得分仅为0.4317。而在引入向量相似度检索、监督微调与强化学习组成的三阶段结构后, 模型得分提升至0.4705, 验证了结构化增强策略在中药推荐任务中的有效性。

方法	Task2-Score
Qwen2.5-7B (直接生成)	0.4318
Qwen2.5-7B (任务1+2联合生成)	0.3869
Qwen2.5-7B (检索+微调+强化学习)	0.4705

表 4: 任务2中药处方推荐方法在验证集上的表现

此外, 本文还探索了将任务1与任务2合并为统一的生成式任务进行训练的方案, 结果表明其性能显著低于本文所采用的“微调+集成”方法。这表明在多任务建模中, 任务结构与标签划分的明确性对模型性能具有关键影响。

总体而言, 本文提出的方法在两个子任务上均显著优于baseline, 验证了所设计的模型架构、训练流程及模块组合在中医智能诊疗场景中的实用性与有效性。

5 总结与展望

本文面向CCL25-Eval 任务9中的中医辨证辨病与中药处方推荐问题, 提出了一套基于大语言模型的系统性解决方案。在子任务1中, 本文采用Qwen2.5-7B、Mistral-7B 和Baichuan2-7B 三种预训练模型, 结合QLoRA 进行高效参数微调, 并通过模型集成投票策略, 实现了多标签的证型与疾病精准预测。在子任务2中, 本文构建了融合向量检索、监督微调与强化学习的中药处方推荐框架, 在有限标注条件下实现了稳定且准确的推荐效果。

尽管取得了较好结果, 本方法仍存在改进空间。首先, 子任务1尚可进一步引入医学知识图谱等结构化知识, 以增强模型的可解释性与对中医语义的显式建模能力; 其次, 子任务2中的推荐机制未显式建模药物配伍原则与禁忌规则, 未来可结合中医理论约束提升处方的合理性与临床适用性。此外, 随着大模型推理能力的持续增强, 如何与中医辨证论治的逻辑深度融合, 成为实现中医人工智能可信化、精准化发展的关键方向。

参考文献

- Zafar Ali, Yi Huang, Irfan Ullah, Junlan Feng, Chao Deng, Nimbeshaho Thierry, Asad Khan, Asim Ullah Jan, Xiaoli Shen, Wu Rui, et al. 2023. Deep learning for medication recommendation: a systematic survey. *Data Intelligence*, 5(2):303–354.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- JaWanna Henry, Yuriy Pylypchuk, Talisha Searcy, Vaishali Patel, et al. 2016. Adoption of electronic health record systems among us non-federal acute care hospitals: 2008–2015. *ONC data brief*, 35(35):2008–15.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2023. Mistral 7b. *arXiv e-prints*, page arXiv:2310.06825.
- Yiqiao Jin, Mohit Chandra, Gaurav Verma, Yibo Hu, Munmun De Choudhury, and Srijan Kumar. 2024. Better to ask in english: Cross-lingual evaluation of large language models for healthcare queries. In *Proceedings of the ACM Web Conference 2024*, pages 2627–2638.
- Enkelejda Kasneci, Kathrin Seifler, Stefan K  chemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan G  nnemann, Eyke H  llermeier, et al. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274.
- Junling Liu, Ziming Wang, Qichen Ye, Dading Chong, Peilin Zhou, and Yining Hua. 2023. Qilin-med-vl: Towards chinese large vision-language model for general healthcare. *arXiv preprint arXiv:2310.17956*.
- Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, Zijian Zhang, Feng Tian, and Yefeng Zheng. 2024. Large language model distilling medication recommendation model. *arXiv preprint arXiv:2402.02803*.
- Junyuan Shang, Cao Xiao, Tengfei Ma, Hongyan Li, and Jimeng Sun. 2019. Gamenet: Graph augmented memory networks for recommending medication combination. In *proceedings of the AAAI Conference on Artificial Intelligence*, pages 1126–1133.
- Chenhao Su, Sheng Gao, and Si Li. 2020. Gate: graph-attention augmented temporal neural network for medication recommendation. *IEEE Access*, 8:125447–125458.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timoth  e Lacroix, Baptiste Rozi  re, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Paul U Unschuld. 2010. *Medicine in China: a history of ideas*, volume 13. Univ of California Press.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Chaoqi Yang, Cao Xiao, Fenglong Ma, Lucas Glass, and Jimeng Sun. 2021. Safedrug: Dual molecular graph encoders for recommending effective and safe drug combinations. *arXiv preprint arXiv:2105.02711*.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. 2023. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2. 5 technical report. *arXiv e-prints*, pages arXiv–2412.

- Songhua Yang, Hanjie Zhao, Senbin Zhu, Guangyu Zhou, Hongfei Xu, Yuxiang Jia, and Hongying Zan. 2024b. 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*, pages 19368–19376.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*.
- Heyi Zhang, Xin Wang, Zhaopeng Meng, Zhe Chen, Pengwei Zhuang, Yongzhe Jia, Dawei Xu, and Wenbin Guo. 2024. Qibo: A large language model for traditional chinese medicine. *arXiv preprint arXiv:2403.16056*.
- Yizhen Zheng, Huan Yee Koh, Maddie Yang, Li Li, Lauren T May, Geoffrey I Webb, Shirui Pan, and George Church. 2024. Large language models in drug discovery and development: From disease mechanisms to clinical trials. *arXiv preprint arXiv:2409.04481*.
- Xian Zhou, Sai Wang Seto, Dennis Chang, Hosen Kiat, Valentina Razmovski-Naumovski, Kelvin Chan, and Alan Bensoussan. 2016. Synergistic effects of chinese herbal medicine: a comprehensive review of methodology and current research. *Frontiers in pharmacology*, 7:201.