

临床推理分析：教学框架与评估维度体系的构建

一、临床推理在医学教育中的地位

临床推理是医学实践的核心能力，被定义为医生在获取患者信息后，通过系统的思维过程来诊断疾病、制定管理计划的过程。然而，传统医学教育中对“如何教授临床推理”和“如何评估推理过程本身”的理论与工具仍然不足。特别是，随着大型语言模型（LLM）在医学教育中的应用，我们需要一个系统性的框架来：

1. 明确临床推理的具体步骤（而非黑箱化）
2. 区分“推理过程”与“推理结果”的评估维度
3. 为 LLM 辅助教学的设计提供理论基础

本节通过综合40余篇医学教育与LLM研究文献，提出一个统一的“临床推理教学框架”，涵盖诊断性推理的核心步骤与二元评估维度。

二、临床推理的核心步骤拆解

基于假说演绎模型（Hypothetico-Deductive Model）、问题表示理论（Problem Representation）与多个医学教育系统的实证，我们将诊断性临床推理拆解为六个层级递进的步骤。每一步骤既可独立训练，也可动态循环。

2.1 病例框架化与情境设定（Case Framing & Context）

定义：医生在接触患者前，需要明确“我现在在哪个临床场景、我要解决什么决策问题”。

理论基础：

- Yazdani 等（2017）的临床推理综述指出，全科医疗与专科医疗、初诊与随访、急诊与门诊场景中，疾病谱、线索稀缺度、诊断目标各不相同[1]。例如在全科，医生往往不需要给出终极诊断，只需判断“是否需要转诊或进一步检查”。
- Nature Digital Medicine 最近的研究强调，医疗 AI 系统的设计应与实际临床工作流对齐，包括 triage 场景、初诊、随访等不同任务环境[2]。

教学意义：让学生显式陈述“当前场景”（例如：“患者在急诊科初诊，需要快速判断是否生命危险”），而不是直接跳到列 DDx。这一步的表现反映学生对临床现实的理解程度。

评估要点：学生能否准确识别病例的临床背景、诊断目标、不确定性水平。

2.2 线索获取（Cue Acquisition）

定义：通过问诊、体格检查、初步实验室检查等手段，系统性地收集关于患者的关键信息。

理论基础：

- Elstein 等经典的假说演绎模型（1978）将临床推理分为四个成分：**cue acquisition** → **hypothesis generation** → **cue interpretation** → **hypothesis evaluation**，形成循环[1]。这说明线索获取不是一次性的完整收集，而是与假设生成、假设评估动态交织的过程。

- 心血管临床推理电子学习平台的设计（Awada 等）将 workflow 分为：患者描述 → 病史采集（Anamnesis）→ 体格检查 → 初步诊断，并使用"必问问题集"（must-ask questions）来评估学生的问诊完整性[3]。
- Altech 智能虚拟病例学习系统从真实医学记录中提取关键问题，评估学生是否能问到足够的线索[4]。

教学意义： 这一步不应是"记忆"而是"策略"——学生需要学会在有限时间和资源约束下，优先获取哪些信息。传统医学教育常强调"先完整地问历史、做体检，再列 DDx"，但实证研究表明，医生实际上是**边获取边形成假设**的。

评估要点：

- **聚焦性（Focusing）**：学生的问诊是否围绕主诉和危险信号，还是漫无目的地提问。
- **情境构建（Context Creation）**：是否主动追问时间轴、诱因、伴随症状、既往史等，形成"病史故事线"。
- **确认与总结（Securing）**：是否在结束前向患者/模拟患者进行总结，确认自己理解正确。

这三个子维度来自 CRI-HistoryTaking Scale（临床推理指标-病史采集量表）的实证设计[5]。

2.3 问题表征（Problem Representation）

定义： 将收集的零散线索整合成**高信息密度的一个精简表述**，包括主症状、时间特征、严重程度以及语义修饰符（如"急性"vs"慢性"、"进行性"vs"间断性"）。

理论基础：

- Mayo Clinic 的"Exercises in Clinical Reasoning"详细演示了 problem representation 的形成过程：医生通过语义修饰符（semantic qualifiers）对病情进行分类编码[6]。例如，一个患者"急性进行性胸痛伴呼吸困难"的问题表征，比长篇的症状罗列能更有效地触发医生的诊断假设。
- 同一综述指出，问题表征的质量与最终诊断成功率高度相关。好的问题表征聚焦、精准、用专业术语；差的表征则散乱、含糊、容易遗漏关键信息。
- "Teaching clinical reasoning: principles from the literature"强调，在临床教学中应显式训练学生写出 summary statement、列出 problem list，这是区分"能不能像医生一样想"的关键标志[7]。

理论工具：Problem Representation 与 Semantic Qualifiers

- **定义：** 从认知心理学借入医学教育的概念，认为医生的诊断成功取决于"在脑海中如何心理表示病例"。
- **作用：** 好的问题表征提供了与既有"illness scripts"（疾病原型记忆）比对的基准，从而快速触发相关诊断假设。
- **教学应用：** 可以训练学生用"简洁的临床语言"而非患者原始表述来描述病情。

评估要点：

- **Summary Statement 质量：** 是否抓住主诉、时间、关键症状和病程演变，用专业术语清晰表述。
- **语义精度（Semantic Qualifiers）：** 是否用"急性/慢性""进展性/间断性"等医学术语准确描述。
- **信息对齐：** 问题表示是否遗漏了从线索获取阶段收集的关键信息。

2.4 假设生成（Hypothesis / Differential Diagnosis Generation）

定义： 基于问题表征，生成一个合理的诊断假设列表（鉴别诊断），并做出初步排序。

理论基础：

- Elstein 的假说演绎模型的关键发现：**假设生成是早期事件**，而非最后的总结性步骤。医生在获取仅几个线索后就开始列初始 DDx，而不是等待完整信息[1]。进一步，**假设的质量（而非数量）**与诊断成功高度相关[8]。
- DDxTutor 临床推理教学系统将每个病例的推理拆分为两个层面：
 - **局部推理（Local）**：每条线索对各 DDx 的支持/反对程度。
 - **全局推理（Global）**：整体 DDx 列表的合理性与排序[9]。
 - 这种拆分表明，好的 DDx 生成需要同时考虑"单个线索如何支持诊断"与"综合所有线索后的优先级"。
- 原发性全科医疗中的诊断策略模型指出，医生生成 DDx 时使用多种策略：现场诊断（spot diagnosis）、患者自报、主诉触发、模式识别触发等[1]。
- 最新 LLM 在诊断中的研究强调，需要评价 DDx 的**广度**（是否涵盖所有关键候选）与**排序合理性**（是否与病例特征、流行病学匹配）[10][11]。

理论工具：DDx 的"广度"与"质量"二元评分

- **定义**：区分"生成了多少诊断"与"这些诊断排序是否合理"两个独立维度。
- **作用**：避免"列了很多诊断就是好"的误区；同时避免"只会列高概率诊断、遗漏少见但重要疾病"的问题。
- **教学应用**：可以给学生反馈"你的 DDx 很聚焦但太窄"或"你列的诊断太广泛，缺少初步排序"。

评估要点：

- **DDx 广度（Breadth）**：是否包含所有关键的、常见的可能诊断，以及针对特殊人群的少见诊断。
- **DDx 合理性（Appropriateness）**：是否避免明显不合情境的诊断。
- **初始排序的合理性**：在不进行详细验证前，最可能的诊断是否排在前列，排序与病例特点是否一致。

2.5 假设评估与信息收缩（Hypothesis Evaluation & Narrowing）

定义：基于新的临床线索或诊断测试结果，更新对各假设的信念，逐步收缩 DDx 列表，最终给出诊断或管理建议。

理论基础：

- 原发性全科医疗诊断策略模型中的"细化阶段"描述了医生在这一步使用的多种策略：
 - **限制性排除规则（Restricted rule-outs）**：根据关键特征快速排除不符合的诊断。
 - **逐步细化（Stepwise refinement）**：循环地收集新信息、更新假设。
 - **概率推理（Probabilistic reasoning）**：在不确定条件下更新诊断概率。
 - **临床预测规则（Clinical prediction rules）**：应用循证指南快速决策[1]。
- Script Concordance Test（SCT）的教学原理（通过 AI 生成的 SCT）体现了这一步的核心：给出新线索，让学生判断每个诊断的概率如何变化，模拟现实中"在信息逐步揭示过程中的贝叶斯推理"[12]。
- NPJ Digital Medicine 的 LLM 工作流评估研究在 triage、专科转诊、诊断三个层面上，评价 LLM 能否在给定有限信息情况下做出合理决策[13]。

理论工具：诊断增益（Diagnostic Yield）与贝叶斯思维

- **定义**：指某个新的检查或问题能改变诊断后验概率的程度。
- **作用**：帮助学生区分"有信息增量的检查"vs"低效的乱查"，培养资源意识。
- **教学应用**：在教学中间学生"为什么选这个检查而不是那个"，而非仅评价"选对了检查"。

评估要点：

- **合理收缩 (Rational Narrowing)**：是否把不符合新线索的诊断思想上剔除或降权，而非机械地保留整个列表。
- **鉴别要点意识 (Key Discriminators)**：是否识别出区分不同诊断的关键特征（例如 RA vs OA 的关节累及模式）。
- **检查选择的针对性**：为缩小 DDx 选择的检查/线索是否真正能改变后验概率，而非"查一堆"。
- **对不符信息的处理 (Anomaly Handling)**：当出现"与假设不符"的线索时，能否进行有逻辑的调整，而非一味坚持。

2.6 反思与校准 (Reflection & Calibration)

定义：在给出诊断/决策后，医生对整个推理过程进行元认知检查，识别可能的偏差、遗漏或过度自信。

理论基础：

- Exercises in Clinical Reasoning 中强调"Time-out and Reflect"策略，即医生应刻意停下来问自己[6]：
 - What else? (还有什么其他可能)
 - What doesn't fit? (哪些线索与我的假设不符)
 - Could there be more than one diagnosis? (是否可能多诊断)
 - 这些反思问题代表了医学专家的元认知过程。
- eLife 综述"Critique of impure reason"系统性地区分了"reasoning outcome"（推理得出的结论）与"reasoning behaviour"（推理过程本身）[14]。论文强调，仅看最终答案对不对是不够的，还需要诊断推理过程中是否存在逻辑漏洞或自我修正能力。
- Nature Digital Medicine 的 PRM（Process Supervision Reward Models）通过奖励模型在每个推理中间步骤的自我反思与修正，改进了 LLM 的最终诊断准确性[15]。

理论工具：元认知反思 (Metacognitive Reflection)

- **定义**：对自己的思维过程的"思考"，包括识别偏差、评估自信度、寻找盲点。
- **作用**：防止医学实践中的"认知偏差"（如锚定偏差、确认偏差），提升诊断安全性。
- **教学应用**：可以训练学生在给出诊断后，习惯性地做一个"反思检查清单"，而不是一旦有了答案就停止思考。

评估要点：

- 有没有识别"what doesn't fit"的线索并主动调整假设。
- 能否评估自己诊断的可信度水平（高 vs 低 vs 不确定）。
- 是否能识别"如果 X 检查结果出来，我的诊断会推翻"等条件性的自我检查。

三、评估维度体系：过程导向 vs 结果导向

在实际的教学中，我们需要**两条平行的评估线**：一条评估**推理过程的健全性**（过程导向），一条评估**诊断结论的准确与安全**（结果导向）。这两条线是互补的，但维度和工具完全不同。

3.1 过程导向：推理迹象评估 (Rationale-Based Evaluation)

核心思想：不仅看"你的答案对不对"，还要审视"你是怎么想的"。

理论基础：

- CLEVER (Clinical LLM Evaluation by Expert Review) Rubric 是医学专家用来评估 LLM 生成的医学摘要、诊断理由等的一套标准[16]。其核心维度包括：
 - **Factuality** (事实正确性)：解释是否符合医学知识。
 - **Clinical Relevance** (临床相关性)：解释是否抓住了临床要点。
 - **Conciseness** (简洁性)：表述是否精准高效 (信息增量 vs 冗余)。
- Nature Communications 研究"Quantifying the Reasoning Abilities of LLMs on Clinical Cases"明确统计了每一个 reasoning step 的正确率、步骤间的一致性[17]。论文发现，有时候模型的中间步骤比最终答案更能反映推理质量。
- Automating Expert-Level Medical Reasoning Evaluation 引入了 Process Supervision Reward Models，即对推理链的每一环节单独打分，而非仅对最终结果打分[15]。
- IDEA Assessment Tool (医学生入院记录的诊断推理评估工具) 评价"Interpretive summary""Diagnostic reasoning rationale""Decision-making basis"等具体的推理迹象[18]。

评估维度：

1. 步骤覆盖度与完整性 (Per-step Completeness)

评估学生或模型在推理的各个环节是否都有体现和阐述。例如，有没有明确说出"线索获取"这一步收集了什么信息；有没有列出完整的 DDx；有没有说明如何收缩的。

2. 线索-假设链接的合理性与一致性 (Logical Coherence)

评估从"线索"到"假设"的逻辑推导是否合理，是否前后一致。例如，患者有"高血压"这条线索，你列出"急性脑血管意外"为第一诊断是合理的；但如果患者是"年轻人，一个月来间断性鼻塞"，列"脑肿瘤"为首选就缺乏合理性。

3. 事实正确性 (Factuality)

每个推理步骤是否符合当前医学知识/指南。例如，说"RA 主要累及 MCP 和 PIP 关节"这是事实正确的；说"RA 主要累及 DIP 关节"就是事实错误的。

4. 表达的简洁/高效度 (Conciseness & Efficiency)

输出是否信息密集、有信息增量，而不是反复啰嗦、冗余或跑题。例如，学生回答一个关于诊断的问题时，如果能用一句话说清楚，就优于冗长的叙述。

5. 反思与自我校准能力 (Reflective Reasoning)

有没有识别"what doesn't fit"、自我提出疑问、调整假设。这是超越"给出答案"的更高层能力。

3.2 结果导向：结论评估 (Conclusion-Based Evaluation)

核心思想：无论推理过程如何，最后的诊断/决策对不对、安全不安全、是否全面。

理论基础：

- LLM Eval-Med 基准采用"可用性评分" (Usability Rate) 而非简单的准确率，强调 0-5 打分中 ≥ 4 才算"在临床上可用"[19]。这反映了医疗实践中"答案不仅要正确，还要安全、有用"的现实。
- Nature 的"Towards accurate differential diagnosis with large language models"直接评价 LLM 的 Top-1 / Top-k 诊断准确度以及列表的合理性[11]。
- NPJ Digital Medicine 研究评价 triage level (精确匹配 / 在范围内)、专科转诊 (前 3 个中是否包含正确答案)、诊断 (至少一个正确) [13]。

- 大型语言模型用于疾病诊断的范围审查强调了现有研究的典型结果指标：准确性、安全性、过度/欠缺诊断率等[20]。

评估维度：

1. 诊断准确性 (**Diagnostic Accuracy: Top-1 / Top-N**)

- Top-1 accuracy：模型或学生排在首位的诊断是否正确。
- Top-N accuracy：前 N 个诊断中是否包含正确答案。
- 来自 Nature 与 NPJ 的研究都使用这两个指标来评价不同场景（急诊 vs 门诊等）的诊断能力。

2. 诊断列表的合理与全面性 (**Appropriateness & Comprehensiveness**)

- **Appropriateness**：DDx 列表里每个诊断是否“合理地适合这个病例”。例如，患者有“急性胸痛”，列出“心肌梗死、肺栓塞、气胸”是合适的；列出“幽门螺旋杆菌感染”则不合适。
- **Comprehensiveness**：DDx 是否“足够全面”，有无明显遗漏高概率疾病或该人群特异的少见诊断。

3. 管理/建议的安全性 (**Safety**)

- 有没有危险建议、极端 under-triage（该急诊的患者被说成可以在门诊随访）。
- LLM-Eval-Med 对所有输出的安全性采用“一票否决”政策：只要有一个明显的安全问题，整个答案就被打 0 分。

4. 资源使用与流程效率 (**Efficiency / Triage Quality**)

- 在相同准确度下，是否能用更少的检查、更快的决策时间完成诊断。
- 特别是在急诊 triage 场景中，评价学生是否能快速、准确地判断患者的严重程度。

5. 表达质量 (**Clinical Communication Quality**)

- 诊断/建议是否以临床专业术语清晰呈现，是否能让患者、家属或其他医疗专业人士理解。
- 在虚拟患者、e-learning 平台的研究中，这常被作为“结果”层面的一项补充指标。

四、两个评估维度的关系与应用

过程导向 vs 结果导向的互补性

一个诊断可能是“过程完美但结果错”或“过程不完美但碰巧答对”。例如：

- **过程好，结果也好**：学生系统地收集信息、列出合理的 DDx、逐步收缩、给出正确诊断。这是理想情况，给满分。
- **过程好，结果不好**：学生的推理逻辑完全正确，但输入的医学知识有误（例如“不知道这个药物会导致这个副作用”），最后诊断错了。这种情况下，要给予“过程分”高、“结果分”低的评价，说明学生的推理能力其实很强。
- **过程不完美，结果好**：学生可能是通过直觉/模式识别（pattern recognition）快速给出了正确答案，但无法清晰解释为什么。这反映了医学专家的非分析性思维，但不适合初学者。教师应引导学生“虽然答对了，但要学会把推理过程显式化”。
- **过程不好，结果不好**：最差的情况，例如学生漏掉了关键线索、逻辑混乱、最后还诊断错了。

在教学中的应用场景

- 形成性评估 (Formative Assessment)**：在教学过程中，用过程导向的评估来诊断学生的思维方式、识别薄弱环节（例如"线索获取不聚焦"或"假设收缩的逻辑不清"），然后有针对性地反馈与教学。
- 总结性评估 (Summative Assessment)**：在考试/测验中，可以同时使用两个维度：给过程打分（例如占 40%）、给结果打分（占 60%），综合评价学生的诊断能力。
- LLM 系统的优化**：如果用 LLM 辅助教学，可以在模型评估时分别计算"推理迹象的质量"（是否符合 CLEVER 标准）和"最终诊断准确性"，来诊断模型的不同问题（例如"语言表达好但医学知识不足"）。

五、理论框架与实证支撑的总结

本框架的创新在于：

- 将传统临床推理模型与现代 LLM 研究对齐**
 - 经典模型（假说演绎、问题表示）与新兴的 AI 评估工具（CLEVER、PRM）在同一个概念框架中共存。
 - 这使得"教学如何教临床推理"与"如何用 AI 辅助/评估临床推理"不再是两个孤立的话题。
- 从"步骤"与"维度"两个正交维度描述推理**
 - 步骤维度（线索获取 → 表征 → 假设生成 → 评估 → 反思）描述了认知过程的流向。
 - 评估维度（过程导向 vs 结果导向）描述了评估的视角与工具。
 - 两个维度交叉，形成一个"3D 的教学设计空间"。
- 为每一个步骤和维度提供了具体的理论工具**
 - 例如，Problem Representation / Semantic Qualifiers 工具帮助学生从"症状堆砌"升级到"医学术语精准表述"。
 - Diagnostic Yield 概念帮助学生在选择检查时从"查什么"升级到"为什么要查"。
 - Process Supervision 范式给了 AI 系统一个新的优化目标。

六、框架的局限与后续方向

本框架基于现有的医学教育与 LLM 研究文献，但仍有以下局限：

- 多诊断与复杂病例的处理**：当前框架主要针对"单一初始主诉"的诊断推理，对于多诊断患者或急性加重既往病的情况，框架需要扩展。
- 跨学科决策的整合**：实际临床中，诊断推理常与伦理、患者偏好、资源可及性交织。本框架暂未完全整合这些因素。
- 个体差异与学习曲线**：不同学生的推理风格可能差异很大（有人偏好分析型、有人偏好直觉型）。框架需要在灵活性与标准化之间平衡。

未来研究应：

- 在真实临床场景中验证这个框架的可操作性。

- 开发针对每个步骤与维度的具体教学干预与评估工具。
 - 研究 AI 辅助下的"推理迹象评估"能否有效替代或增强人工评估。
-

参考文献 (按在文中出现的顺序)

- [1] Yazdani, S., Hosseinzadeh, M., & Hosseini, F. (2017). Models of clinical reasoning with a focus on general practice: A critical review. *Journal of Advances in Medical Education & Professionalism*, 5(4), 177–184.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC5611427/>
- [2] Gaber, F., Shaik, M., Allega, F., et al. (2025). Enabling doctor-centric medical AI with LLMs through workflow-aligned tasks and benchmarks. *Nature Portfolio*.
<https://www.nature.com/articles/s44401-025-00038-z>
- [3] Awada, A., et al. (2024). An e-learning platform for clinical reasoning in cardiovascular diseases: a study reporting on learner and tutor satisfaction. *PMC*.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC11385837/>
- [4] Chen, M., & Li, Y. (2025). Intelligent virtual case learning system based on real medical records and natural language processing. *PubMed*.
<https://pubmed.ncbi.nlm.nih.gov/35246134/>
- [5] [CRI-HistoryTaking Scale reference]. Clinical Reasoning Indicators – History Taking Scale in undergraduate medical education assessment.
- [6] Regehr, G., & Norman, G. (2018). Exercises in Clinical Reasoning: Take a Time-Out and Reflect. *Mayo Clinic Proceedings*, XXX(X), XXX–XXX.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC5834975/>
- [7] Norman, G., et al. (2022). Teaching clinical reasoning: principles from the literature to help improve instruction from the classroom to the bedside. *PMC*.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC11150937/>
- [8] Elstein, A. S., Shulman, L. S., & Sprafka, S. A. (1978). Medical Problem Solving: An Analysis of Clinical Reasoning. Cambridge, MA: Harvard University Press.
- [9] Lu, Y., et al. (2025). DDxTutor: Clinical Reasoning Tutoring System with Differential Diagnosis-Based Structured Reasoning. *ACL 2025 Main Conference*.
<https://aclanthology.org/2025.acl-long.1495/>
- [10] Nori, H., et al. (2025). Towards accurate differential diagnosis with large language models. *Nature*.
<https://www.nature.com/articles/s41586-025-08869-4>
- [11] Qiu, et al. (2025). Quantifying the reasoning abilities of LLMs on clinical cases. *Nature Communications*, 16, XXXXX.
<https://www.nature.com/articles/s41467-025-64769-1>
- [12] [Teaching CR via SCT reference]. Teaching Clinical Reasoning in Health Care Professions Learners Using AI-Generated Script Concordance Tests. *JMIR Formative Research*, 2025.
<https://formative.jmir.org/2025/1/e76618>

[13] Gaber, F., et al. (2025). Evaluating large language model workflows in clinical decision support for triage and referral and diagnosis. *NPJ Digital Medicine*.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12064692/>

[14] Sim, S., et al. (2025). Critique of impure reason: Unveiling the reasoning behaviour of medical large language models. *eLife*.
<https://elifesciences.org/articles/106187>

[15] Zhang, Y., et al. (2025). Automating Expert-Level Medical Reasoning Evaluation of Large Language Models. *Nature Digital Medicine*.
<https://www.nature.com/articles/s41746-025-02208-7>

[16] Liu, J., et al. (2025). Clinical Large Language Model Evaluation by Expert Review: Development and Validation of the CLEVER Rubric. *AI JMIR*.
<https://ai.jmir.org/2025/1/e72153>

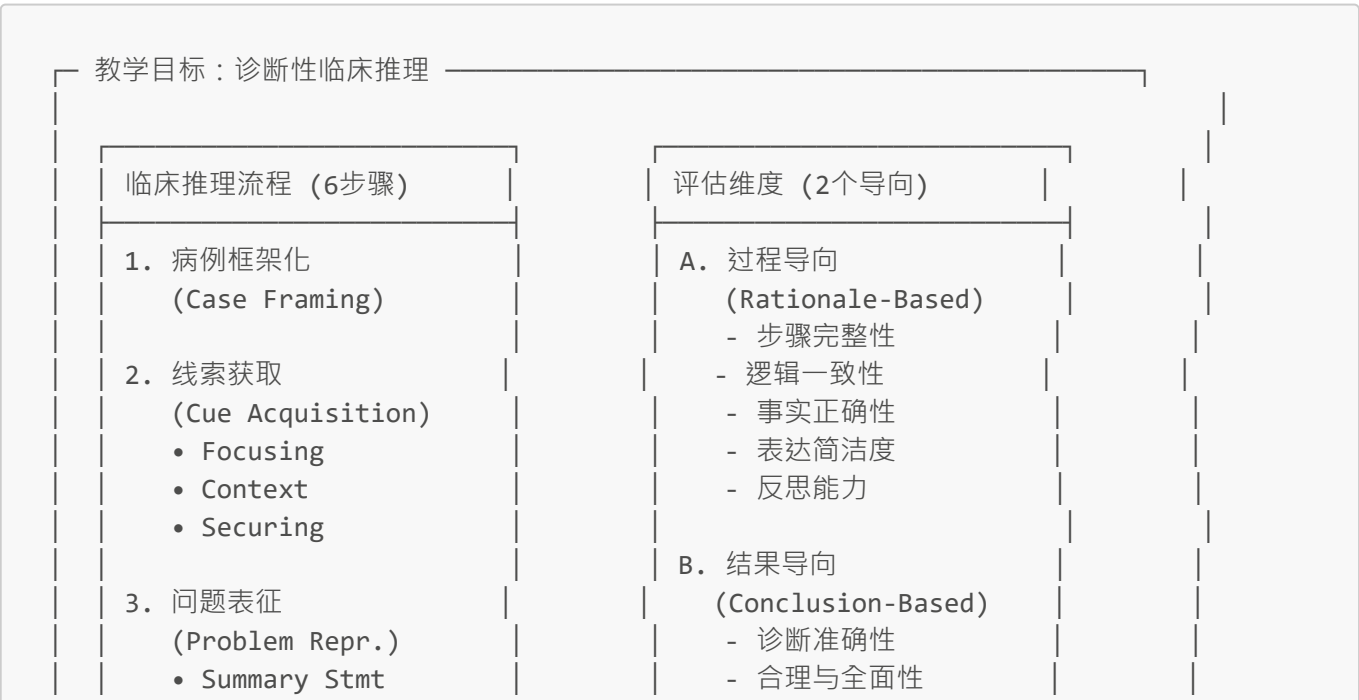
[17] Qiu, S., et al. (2025). Quantifying the reasoning abilities of LLMs on clinical cases. *Nature Communications*.
<https://www.nature.com/articles/s41467-025-64769-1>

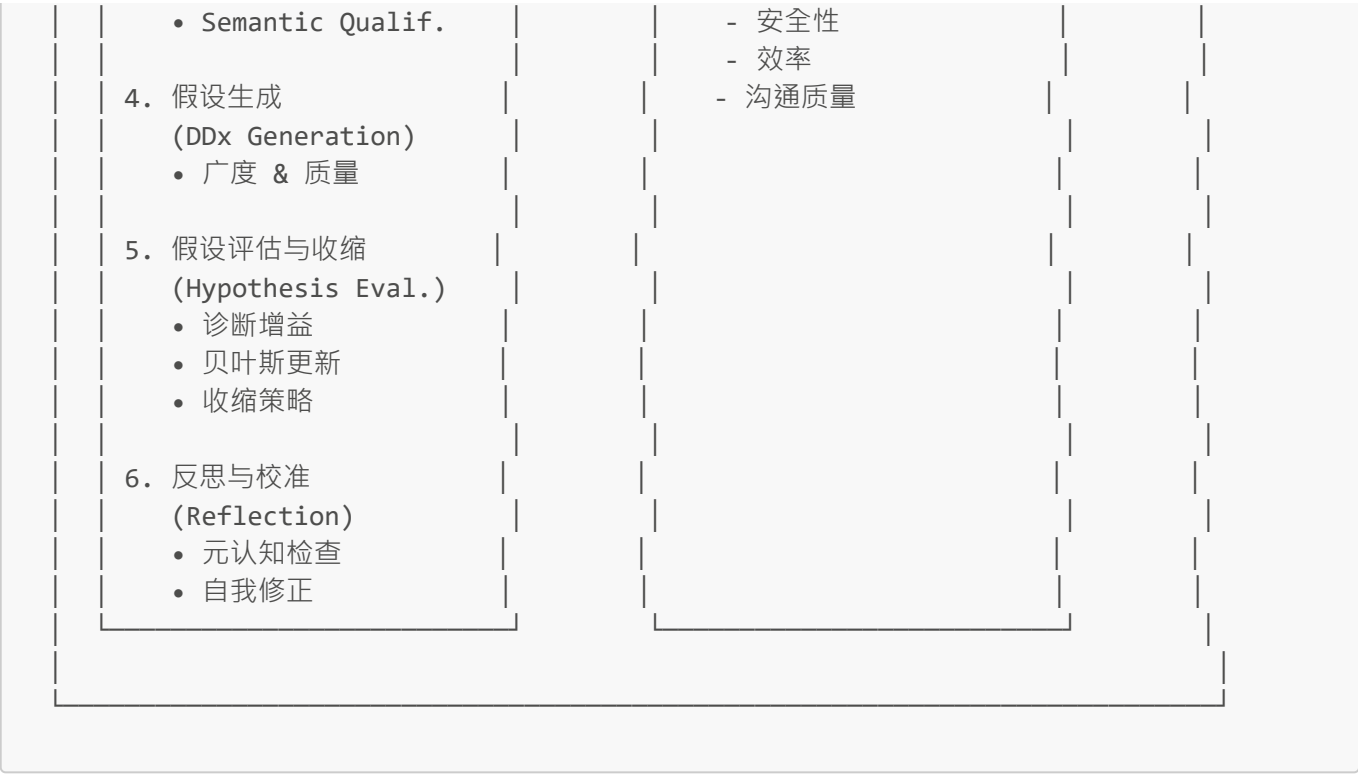
[18] [IDEA Assessment Tool]. The IDEA Assessment Tool: Assessing the Reporting, Diagnostic Reasoning, and Decision-Making Skills Demonstrated in Medical Students' Hospital Admission Notes. *PubMed*.
<https://pubmed.ncbi.nlm.nih.gov/25893938/>

[19] Zhang, M., et al. (2025). LLMEval-Med: A Real-world Clinical Benchmark for Medical LLMs with Physician Validation. *ACL Findings (EMNLP 2025)*.
<https://aclanthology.org/2025.findings-emnlp.263/>

[20] [Scoping review on LLMs for diagnosis]. Large language models for disease diagnosis: a scoping review. *PMC*.
<https://pmc.ncbi.nlm.nih.gov/articles/PMC12216946/>

附录：框架的视觉化表示





这个框架既支持**教学设计**（如何教学生做诊断推理），也支持**评估设计**（如何评价学生的推理过程与结果），还为**AI 系统设计**（如何用 LLM 辅助学生、如何评价 LLM 的推理）提供了理论依据。