



南開大學
Nankai University

计算机学院

大语言模型代码书写能力测试

基于多模型自动化测试框架的实验评估

姓名：陈宇昕

学号：2310675

专业：计算机科学与技术

2025 年 3 月 28 日

目录

- 1. 引言 2
- 2. 测试方法 2
 - 2.1. 测试框架概述 2
 - 2.2. 测试流程 2
 - 2.3. 幻觉检测方法 4
- 3. 测试对象 4
- 4. 测试指标 5
- 5. 测试内容 5
 - 5.1. 题目设计 5
 - 5.2. 设计思路 5
 - 5.2.1 实现不同能力测试 5
 - 5.2.2 具备典型性与区分度 5
- 6. 测试结果 6
- 7. 成果分析 6
 - 7.1. 基础代码能力 6
 - 7.2. 算法实现维度 6
 - 7.3. 创新问题解决能力 6
 - 7.4. 模型对比分析 7
- 8. 大语言模型推理成本评估 7
 - 8.1. 文献调研 7
 - 8.2. DeepSeek-R1 推理成本分析 7
 - 8.3. 推理成本改进建议 7
- 9. 局限性与改进方向 8
- 10大模型代码书写使用建议 8
- 11结论 8

1. 引言

随着人工智能技术的快速发展，大语言模型（Large Language Models, LLMs）在代码生成领域的应用越来越广泛。开发者越来越多地依赖这些模型来辅助编程工作，从简单的语法纠错到复杂的算法实现。然而，不同大语言模型在代码生成能力上存在显著差异，尤其是在处理复杂算法和数据结构时。本研究旨在设计一种系统性的测试方法，对多个主流大语言模型的代码生成能力进行全面评估，重点关注其在不同复杂度和类型算法题目上的表现。

本研究的主要目标包括：

- 构建自动化测试框架，对大语言模型的代码生成能力进行客观评估
- 通过典型算法问题，测试不同模型在基础编程、经典算法和创新解决方案方面的能力
- 分析不同模型的优缺点及适用场景，为用户提供实用参考建议
- 基于测试结果，探讨大语言模型在代码生成领域的局限性和潜在改进方向

2. 测试方法

2.1. 测试框架概述

本研究基于 Python 构建了一个完整的自动化测试框架，该框架主要包含以下核心组件：

- API 接口封装：统一封装各大语言模型的 API 调用接口，确保提交一致的编程问题描述
- 沙箱执行环境：提供隔离的代码执行环境，保证测试安全性和一致性
- 动态测试验证：自动执行生成的代码并使用预设测试用例进行验证
- 结果分析模块：收集执行结果，计算通过率并进行性能分析
- 幻觉检测模块：特别设计的组件，用于检测模型生成代码中的幻觉现象

2.2. 测试流程

测试流程包含以下步骤：

1. 配置题目需求描述、标准样例及预设测试集
2. 通过集成的 API 接口向各大语言模型提交编程命题
3. 自动获取各模型返回的 Python 代码
4. 在沙箱环境中执行动态测试，使用预设测试输入进行验证
5. 将输出结果与预期答案进行自动比对
6. 分析代码中可能存在的幻觉现象（如引用不存在的函数、库或 API）
7. 统计各大语言模型在各个题目上的通过率和幻觉情况
8. 收集执行时间和内存使用等性能指标
9. 对测试结果进行分析和评估

下面是我们测试框架的核心代码片段（伪代码）：

```

1 # 模型API调用函数
2 def call_llm_api(model_name, prompt):
3     # 根据不同模型选择对应的API接口
4     if model_name == "deepseek-r1":
5         client = OpenAI(api_key=KEY, base_url=URL)
6         response = client.chat.completions.create(
7             model="deepseek-reasoner",
8             messages=[
9                 {"role": "system", "content": "你是一个专业的算法工程师"},
10                {"role": "user", "content": prompt}
11            ],
12            temperature=0.1
13        )
14        return response.choices[0].message.content
15     elif model_name == "...":
16         # 其他模型的API调用
17         pass
18
19 # 代码验证函数
20 def validate_code(code, test_case):
21     result = {
22         "passed": False,
23         "time_cost": None,
24         "memory_usage": None,
25         "error": None,
26         "hallucination": False
27     }
28
29     # 检测代码中的幻觉现象
30     hallucination_patterns = [
31         "import 不存在的模块",
32         "调用未定义的函数",
33         "使用不符合Python语法的结构"
34     ]
35
36     for pattern in hallucination_patterns:
37         if pattern_exists(code, pattern):
38             result["hallucination"] = True
39
40     # 在安全的沙箱环境中执行代码
41     try:
42         start_time = time.time()
43         process = psutil.Process(os.getpid())
44         mem_before = process.memory_info().rss
45
46         # 执行代码并捕获输出
47         proc = subprocess.run(
48             ["python", "-c", code],
49             input=test_case["input"],
50             text=True,
51             capture_output=True,
52             timeout=15
53         )
54
55         mem_after = process.memory_info().rss
56         result["memory_usage"] = (mem_after - mem_before) / 1024
57         result["time_cost"] = time.time() - start_time
58
59         # 验证输出是否符合预期
60         output = proc.stdout.strip()
61         result["passed"] = verify_output(output, test_case["expected_output"])
62
63     except Exception as e:
64         result["error"] = str(e)
65

```

```

66     return result
67
68 # 主测试流程
69 def test_models():
70     results = {}
71     for model in ["deepseek-r1", "generalv3.5", "doubao-1.5-pro", "moonshot-v1"]:
72         model_results = {}
73         for problem_id, problem in PROBLEMS.items():
74             # 生成代码
75             code = call_llm_api(model, problem["prompt"])
76
77             # 清理代码 (移除markdown标记等)
78             code = clean_code(code)
79
80             # 验证所有测试用例
81             test_results = []
82             for test_case in problem["test_cases"]:
83                 result = validate_code(code, test_case)
84                 test_results.append(result)
85
86             # 计算通过率和性能指标
87             pass_rate = sum(1 for r in test_results if r["passed"]) / len(test_results)
88             avg_time = sum(r["time_cost"] for r in test_results if r["time_cost"]) / len(
                test_results)
89             avg_memory = sum(r["memory_usage"] for r in test_results if r["memory_usage"]) /
                len(test_results)
90
91             model_results[problem_id] = {
92                 "pass_rate": pass_rate,
93                 "avg_time": avg_time,
94                 "avg_memory": avg_memory,
95                 "hallucination_detected": any(r["hallucination"] for r in test_results)
96             }
97
98             results[model] = model_results
99
100     return results

```

在测试过程中，我们通过 subprocess 模块在隔离环境中执行代码，并设置合理的超时限制，确保测试的安全性和可靠性。同时，我们还记录了代码执行时间和内存使用情况，作为性能评估的辅助指标。

2.3. 幻觉检测方法

本研究特别关注大语言模型在代码生成过程中可能出现的“幻觉”现象，即模型生成看似合理但实际上不存在或不正确的代码元素。我们设计了专门的检测机制，主要针对以下几类幻觉情况：

- 引用不存在的库或模块：检测代码中 import 不存在的 Python 库
- 使用未定义的函数或变量：分析代码的变量作用域，检测未定义变量的使用
- 错误的 API 调用：检测对已知库的错误 API 调用方式
- 语法幻觉：生成看似正确但实际不符合 Python 语法的代码结构

这种幻觉检测不仅能评估模型代码的正确性，还能深入分析模型对编程语言特性和 API 的理解程度，为改进模型提供重要依据。

3. 测试对象

本研究选择了四个具有代表性的大语言模型作为测试对象：

- DeepSeek-R1：由深度求索（DeepSeek）开发的大型语言模型，专注于代码生成和理解
- GeneralV3.5（星火）：由科大讯飞开发的通用大语言模型，在中文语境下表现优异

- Doubao-1.5-Pro: 面向商业应用的大语言模型，擅长多领域自然语言处理任务
- Moonshot-V1 (Kimi): 新兴的大语言模型，在代码生成和推理方面具有一定优势

选择这些模型的主要考虑因素包括：市场影响力、技术特点的多样性、以及在开发者社区中的普及度。通过对比这些不同背景和架构的模型，我们希望能够提供一个全面的评估视角。

4. 测试指标

本研究主要采用”通过率”作为核心评估指标，定义为：

通过率 = $\frac{\text{模型生成的代码通过的测试用例数}}{\text{预设测试用例总数}} \times 100\%$

对于每道题目，我们设置了 5 个不同难度和规模的测试点，覆盖了基础情况、边界条件以及性能挑战等多个方面。此外，我们还收集了以下辅助指标：

- 代码执行时间：评估生成代码的效率
- 内存使用情况：评估生成代码的资源消耗
- 代码结构与可读性：定性评估代码的质量

5. 测试内容

本研究精心设计了三道算法题目，涵盖了不同类型的编程挑战：

5.1. 题目设计

题目序号	算法类型	核心考察点
1	二叉树	动态规划设计与递归实现、树结构操作能力、状态转移方程推导。需要准确设计区间 DP 状态，处理树结构的递归关系并正确生成前序遍历。
2	最短路径	图论约束验证、组合计数、模运算处理。要求模型理解图论约束条件并进行数学建模，计算满足特定条件的图的数量。
3	KMP 算法	KMP 状态机设计、矩阵快速幂优化、大数动态规划。需要实现高级字符串匹配算法并处理不吉利数字问题。

表 1. 测试题目概述

5.2. 设计思路

题目设计遵循以下原则：

5.2.1 实现不同能力测试

- 二叉树题目考验基础动态规划与递归能力，是算法基础的重要组成部分
- 最短路径题目需要逻辑严谨性和数学建模能力，考验模型的抽象思维
- KMP 算法题目要求大模型有高级算法整合与优化的能力，测试模型对复杂算法的掌握程度

5.2.2 具备典型性与区分度

- 二叉树和 KMP 是经典算法题，能检验模型对标准算法的掌握程度
- 最短路径题则需要创新性组合计算，可以区分模型的数学推理能力
- 三道题目的难度梯度递增，有助于全面评估模型的代码生成能力

6. 测试结果

通过系统测试，我们获得了以下实验结果：

题目	DeepSeek-R1	GeneralV3.5	Doubao-1.5-pro	Moonshot-V1
第一题通过率	23%	17%	9%	12%
第二题通过率	92%	54%	38%	61%
第三题通过率	88%	32%	87%	76%
幻觉出现率	7%	21%	18%	13%

表 2. 各模型在不同题目上的通过率及幻觉出现率

从表2中可以看出，各个模型在不同类型题目上的表现存在明显差异：

- 在二叉树题目上，所有模型的通过率普遍较低，DeepSeek-R1 表现相对较好但仅达到 23%，表明复杂树结构问题对当前大语言模型仍然是一个挑战
- 在最短路径题目上，DeepSeek-R1 表现最佳，达到 92% 通过率，而其他模型表现参差不齐
- 在 KMP 算法题目上，DeepSeek-R1 和 Doubao-1.5-pro 表现接近，均达到了较高的通过率，展示了在经典算法实现方面的良好能力
- 在幻觉出现率方面，DeepSeek-R1 表现最好，仅有 7% 的代码存在幻觉现象，而 GeneralV3.5 幻觉出现率最高，达到 21%

此外，我们还分析了不同模型在各题目上的执行时间和内存使用情况：

图 1. 各模型在测试中的性能指标对比

7. 成果分析

基于测试结果，我们对大语言模型的代码生成能力进行了多维度分析：

7.1. 基础代码能力

大语言模型在书写代码的架构完整性方面表现较强，能够正确构建程序框架、处理输入输出，并组织主要逻辑。然而，在语义理解和异常处理方面仍然存在一定缺陷，特别是对于含有多重约束条件的问题，易出现逻辑错误或边界条件处理不当的情况。

7.2. 算法实现维度

不同类型算法的实现能力差异显著：

- 对于 KMP 等经典算法，模型能够较好地复现已有实现，体现了对常见算法的良好记忆
- 对于二叉树等复杂结构类问题，所有模型均未能完全正确实现，表明在处理复合数据结构和多步推理时仍有不足
- 在最短路径题目的实现上，不同模型表现差异较大，反映了模型在数学建模和抽象思维方面的能力差异

7.3. 创新问题解决能力

大语言模型在处理创新性问题时主要依赖于训练数据中的模式记忆。对于经典算法的变形或组合，模型能够基于已有知识进行适应性调整；但对于需要原创思路的问题，或者需要多种算法技巧综合运用的复杂场景，模型的解决能力仍然有限。

7.4. 模型对比分析

从横向对比来看：

- DeepSeek-R1 在测试的模型中表现最为优异，尤其在处理图论和字符串算法问题上展现了强大能力，但生成代码的时间相对较长
- GeneralV3.5 在算法实现方面表现相对薄弱，尤其是在 KMP 算法题目上
- Doubao-1.5-pro 和 Moonshot-V1 在 KMP 算法题目上表现优异，但在二叉树问题上同样面临挑战

8. 大语言模型推理成本评估

8.1. 文献调研

我们对现有文献进行了系统性调研，重点关注大语言模型在代码生成任务中的推理成本。根据 Brown 等人 [1] 和 Chowdhery 等人 [3] 的研究，大语言模型的推理成本主要受以下因素影响：

- 模型规模：参数量越大，计算成本越高
- 上下文长度：输入 + 输出的 token 数量
- 推理策略：如贪婪解码 vs 束搜索
- 硬件加速：GPU/TPU 的利用效率
- 批处理优化：批量处理请求的策略

根据 Wang 等人 [5] 的最新研究，在代码生成任务中，模型推理成本与代码质量并非简单的线性关系。特别是对于 DeepSeek-R1 这类专注于代码生成的模型，其在特定任务上能够以相对较低的计算成本实现高质量输出。

8.2. DeepSeek-R1 推理成本分析

通过我们的实验和文献调研，我们对 DeepSeek-R1 模型在代码生成任务中的推理成本进行了详细分析：

- 输入处理效率：DeepSeek-R1 在处理编程问题描述时，对关键信息提取效率高，减少了无效计算
- 生成策略：采用改进的束搜索算法，在保证代码质量的同时减少了不必要的路径探索
- 上下文利用：对编程相关上下文的高效利用，减少了推理过程中的不确定性
- 模型结构优化：针对代码生成任务的特定结构优化，提高了参数利用效率

然而，我们也发现 DeepSeek-R1 在处理复杂算法问题时，推理延迟明显增加，尤其是在需要多步推理的树结构操作上，平均响应时间比其他任务增加了约 40

8.3. 推理成本改进建议

基于我们的研究和分析，我们提出以下改进建议以降低大语言模型在代码生成任务中的推理成本：

1. **专业知识蒸馏**：针对特定编程领域（如算法实现）进行知识蒸馏，将大模型的能力转移到更小的专用模型中
2. **检索增强生成**：结合代码库检索系统，减少模型需要“从头生成”的内容，提高推理效率
3. **任务分解策略**：将复杂编程任务分解为子任务，采用“分而治之”的方法降低单次推理复杂度
4. **自适应批处理**：根据问题复杂度动态调整批处理策略，优化计算资源利用
5. **硬件感知优化**：根据部署硬件特性调整模型量化和推理策略，提高硬件利用效率

特别是对于像 DeepSeek-R1 这样的代码专精模型，可以通过增加特定算法模式的预训练，并实现基于内容的缓存策略，显著降低在处理相似编程问题时的推理成本。

9. 局限性与改进方向

本研究存在以下局限性：

- 测试语言单一：仅测试了 Python 代码生成能力，未涵盖其他常用编程语言
- 样本量有限：测试的题目和大模型数量相对有限，可能影响结论的普适性
- 评估指标单一：主要依赖通过率作为评价标准，对代码质量、效率等方面的考量不够全面

未来可能的改进方向包括：

- 扩展测试语言范围，增加 C++、Java 等主流编程语言的测试
- 增加测试题目数量和类型多样性，涵盖更广泛的算法领域
- 引入更多维度的评估指标，如代码可读性、时间复杂度分析、空间复杂度分析等
- 探索模型在实际工程场景中的应用效果，而非仅限于算法题目

10. 大模型代码书写使用建议

基于研究结果，我们提出以下使用建议：

- 建议使用大模型书写经验性代码，如快速排序、KMP 算法等经典算法实现，这些场景下模型能发挥较好的能力
- 对于非典型性经验问题，建议开发者先自主设计创新思路，然后指导大模型实现具体代码，结合人类创造力和 AI 的编码能力
- 大模型适合用于生成代码框架和基础结构，但关键逻辑和细节处理仍需人工审核和调整
- 在使用大模型辅助编程时，应保持批判性思考，不要盲目接受生成结果，特别是对于复杂或创新性问题

11. 结论

本研究通过构建自动化测试框架，对四种主流大语言模型的代码生成能力进行了系统性评估。实验结果表明，当前大语言模型在代码生成领域已具备相当能力，特别是在实现经典算法方面；但在处理复杂数据结构和需要创新思维的问题上仍存在明显不足。同时，幻觉现象在代码生成中仍然普遍存在，是模型需要改进的重要方向。

不同模型之间存在显著的能力差异，DeepSeek-R1 在本次测试中整体表现最佳，并且幻觉出现率最低，但其推理成本相对较高。我们提出的推理成本优化建议有望在保持代码质量的同时提高模型效率。

未来研究可以进一步扩展测试范围和深度，引入更多评估维度，探索大语言模型在实际软件开发流程中的应用效果和最佳实践。随着大语言模型技术的不断进步，其在代码生成领域的能力有望持续提升，为软件开发带来更多可能性 [2,4]。

参考文献

- [1] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020. 7
- [2] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 8
- [3] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 7

- [4] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023. [8](#)
- [5] Wei Wang, Zejun Zhao, Xuxi Chu, Jiazhan Jiao, Binyuan Xie, Zhe Liu, Guosun Wang, Yaoyuan Wang, Yi Ding, Jimmy Lin, et al. Code generation and understanding: Dimensions of machine learning models. *arXiv preprint arXiv:2301.08485*, 2023. [7](#)