

# 基于大语言模型的思维链推理

黄倪远；谢绎含；单歆芮；陈逸晗；庞瑛霏

## 1. Introduction

思维链推理（Chain of Thought, CoT）是一个与人工智能领域的自然语言处理技术相关的概念，尤其是在解决复杂问题或执行任务的上下文中。CoT通常指的是一系列逻辑步骤，表现为在达到最终答案或解决方案之前，明确地阐述问题解决过程中的每一步骤。

在某些高级的AI模型，如OpenAI的ChatGPT或Google的LaMDA等，中采用CoT来解决更复杂的问题，例如数学问题、逻辑谜题或要求理解上下文的问题。这些模型不仅给出答案，而且能够明确地展示它们是如何通过一连串的思维步骤来推理得出这个答案的。通过这种方式，AI提供了一条清晰的思维路径，这有助于理解AI的决策过程，以及捕获和纠正在这一过程中可能出现的错误。

CoT作为人工智能利于的一个研究方向和实践方法已经展现出了潜力，尤其在处理复杂任务和问题解答的环境中，具体体现在：

- 增强问题解决能力：**CoT在多步骤推理和解决需要连锁推理的复杂问题中表现出色，如数学问题、物理问题、编程挑战等。
- 提高解释性：**通过阐述解决问题的每个步骤，CoT增强了AI的可解释性和透明度，使得用户理解AI决策过程变得更加容易。
- 错误检测和修正：**CoT提供的步骤可以帮助识别和修正推理过程中出现的逻辑错误或计算错误。

尽管CoT能够提供透明的推理链条，但这并不能保证推理的结论总是正确的。AI仍有可能在某些步骤中产生错误，进而引导到错误的结论。阻碍CoT准确性的因素可能包括但不限于：

- 不完整或错误的知识库：**AI模型依赖于它们训练时接触的数据。如果这些数据中含有不正确或不完整的信息，它可能在推理过程中得出错误的结论。

2. **长推理链的问题累积**：在处理长推理链时，每一步的小错误都可能累积起来导致最终结论的显著偏差。
3. **逻辑错误**：在构建推理链的过程中，AI可能会犯逻辑上的错误，比如错误地应用假设、规则或定律，或者对因果关系做出错误的推断。

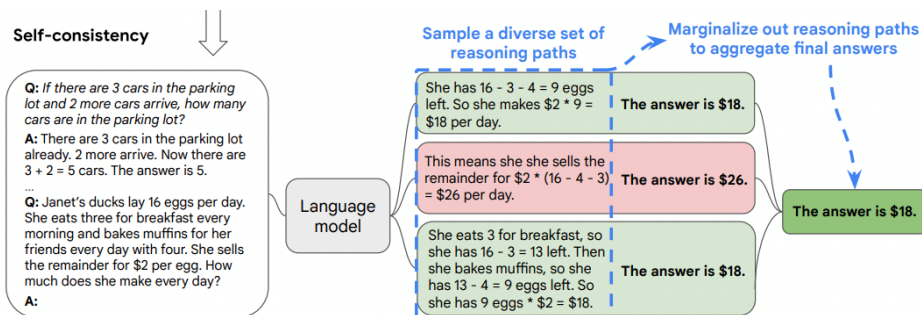
针对以上挑战，我们参照IRCoT以及DIVERSE两种方法对以上缺陷做出改进优化。

## 2. Algorithm

### 2.1 Multi LLM consistency

Multi LLM consistency算法通过整合跨多个大型语言模型（LLMs）的推理路径来提高输出结果的一致性和准确性。该算法的主要步骤如下：

1. 在多个不同的大型语言模型（如GPT4、PaLM2等）上运行相同的思考链提示，以采样出不同的推理路径。这一步骤旨在从多样化的角度探索问题，以便捕捉到可能的不同解决方案和思维方式。
2. 确立共识（quorum），选择最佳响应。在此示例中，使用GPT4作为共识评估器，但也可以采用其他方法。共识的建立是通过比对不同模型给出的答案，从中选出最为一致或出现频率最高的答案，以此作为最终的输出。



- 伪代码如下：

Input: question\_prompt

Output: best\_response (based on consensus)

Function sample\_reasoning\_paths(question\_prompt):

    reasoning\_paths = []

    For each LLM in [GPT4, PaLM2, ...]: # List of predetermined LLMs

```

    response = LLM.generate(question_prompt) # Generate an answer based on the
prompt for each model
    reasoning_paths.append(response)
Return reasoning_paths

Function evaluate_quorum(reasoning_paths):
    Count the frequency of each reasoning path
    Sort reasoning paths by frequency in descending order
    Return the reasoning path with the highest frequency

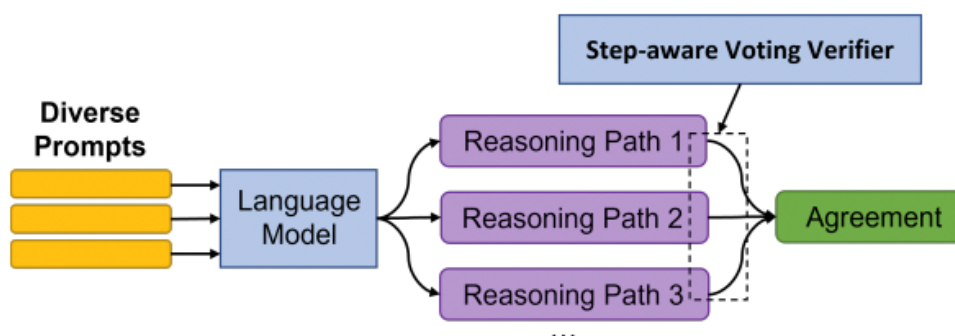
Function main():
    reasoning_paths = sample_reasoning_paths(question_prompt)
    best_response = evaluate_quorum(reasoning_paths)
    Print best_response

# Execute the main function to perform the algorithm
main()

```

## 2.2 DIVERSE

DIVERSE (Diverse Verifier on Reasoning Step) 是一种新颖的方法，旨在增强语言模型的推理能力。该算法主要包括三个组成部分：生成多样化提示以探索不同的推理路径、使用验证器根据加权投票方案过滤不正确的答案、以及逐个验证每个推理步骤而不是整个链条。通过这些组件，DIVERSE能够在推理任务中引导语言模型，提高其性能并取得新的最先进结果。



- 伪代码如下

```


```

```
function DIVERSE(input_question):
    diverse_prompts = generate_diverse_prompts(input_question)
    reasoning_paths = []

    for prompt in diverse_prompts:
        reasoning_path = language_model(prompt)
        reasoning_paths.append(reasoning_path)

    final_answers = voting_verifier(reasoning_paths)

    return final_answers

function generate_diverse_prompts(input_question):
    M1 = select_random_prompts(input_question)
    M2 = sample_reasoning_paths(M1)
    return M2

function select_random_prompts(input_question):
    // 从语言模型中随机选择一组提示
    return random_prompts

function sample_reasoning_paths(prompts):
    reasoning_paths = []
    for prompt in prompts:
        reasoning_path = language_model(prompt)
        reasoning_paths.append(reasoning_path)
    return reasoning_paths

function voting_verifier(reasoning_paths):
    verifier_predictions = []

    for path in reasoning_paths:
        prediction = verifier(path)
        verifier_predictions.append(prediction)

    final_answer = weighted_voting(verifier_predictions)

    return final_answer
```

```

function verifier(reasoning_path):
    // 使用验证器模型对推理路径进行验证
    return prediction

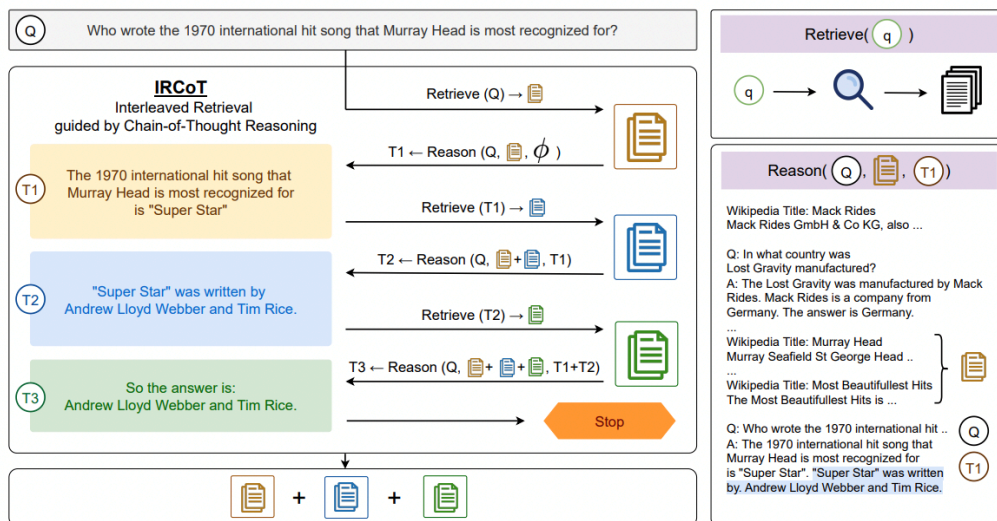
function weighted_voting(predictions):
    // 根据验证器的预测结果进行加权投票
    return final_answer

```

## 2.3 IRCot: 与外部知识结合，利用检索增强技术赋能COT

IRCot, 即 Interleaving Retrieval with Chain-of-Thought Reasoning, 是一个基于结合检索 (Retrieval) 与链式推理 (Chain-of-Thought Reasoning) 的框架, 它专为解决知识密集型的多步问题设计。这种方法通过在检索外部知识与执行链式推理之间交错执行, 旨在提高人工智能系统解决复杂问题的能力。

在传统的链式推理过程中, 模型通常依赖于内部知识库或之前训练中获取的知识来解答问题。然而, 当面临知识密集型的问题时, 内部知识库可能不够全面或更新不及时, 从而限制了推理的准确性和可靠性。



- 算法的步骤概括如下
  1. **问题解析:** 首先对原始问题进行解析, 确定需要回答的核心问题和相关的子问题。

2. **初步推理**: 开始进行CoT推理, 根据已知的内部知识构架答案。此时, 算法可能会识别出需要额外信息以便更好地回答问题的点。
3. **知识检索**: AI系统会在检测到需要更多信息的时候, 自动从指定的外部数据源检索相关知识。这些数据源可能是互联网、科学数据库、知识图谱或者专门的文献库。
4. **信息整合**: 检索到的信息会被整合并结合到连续的思维链中。对于每一步检索到的新信息, 模型都可能需要重新评估先前的结论并进行调整。
5. **反馈循环**: 如果新的信息提示需要进一步的检索或修正推理过程, 算法可能会进入一个反馈循环, 继续交互式地执行检索和推理步骤。
6. **生成回答**: 经过一系列的信息检索和链式推理后, 算法最终输出一个综合了所有必要信息的全面答案。

• 伪代码如下:

```
函数 IRCOT(problem):  
    // 初始化参数, 包括问题、推理状态和动态的知识库  
    current_state ← 初始状态  
    knowledge_base ← 初始知识  
  
    // 循环执行推理步骤, 直到找到满意的答案或达到预定的迭代次数  
    while not finished:  
        // 执行链式推理的当前步骤  
        current_step_answer ← ChainOfThoughtReasoning(problem,  
current_state, knowledge_base)  
  
        // 分析当前推理是否需要外部知识检索  
        if 需要外部知识(current_step_answer):  
            // 确定需要检索的关键词或查询  
            retrieval_query ← 确定检索需求(current_step_answer)  
  
            // 从外部数据源检索知识  
            retrieved_knowledge ← RetrieveKnowledge(retrieval_query)  
  
            // 将检索到的知识加入到动态知识库中  
            knowledge_base ← 更新知识库(knowledge_base,  
retrieved_knowledge)
```

```
// 根据新知识更新当前状态
current_state ← 更新状态(current_state, current_step_answer,
retrieved_knowledge)

// 检查是否已经得出足够的答案或达到终止条件
if 满足结束条件(current_state, knowledge_base):
    finished ← True

// 根据当前状态和知识库生成最终答案
final_answer ← 生成最终答案(current_state, knowledge_base)
返回 final_answer

// 链式推理逻辑
函数 ChainOfThoughtReasoning(problem, state, knowledge):
    // 根据问题、当前状态和知识库推断下一步
    // 返回推理和决定是否需要外部检索的答案

// 确定是否需要外部知识
函数 需要外部知识(step_answer):
    // 分析推理回答，决定是否需要检索
    // 返回布尔值

// 确定检索需求
函数 确定检索需求(step_answer):
    // 分析推理回答，提取出检索关键词或查询
    // 返回检索请求

// 从外部数据源检索知识
函数 RetrieveKnowledge(query):
    // 根据查询从数据源获取信息
    // 返回检索到的知识

// 更新知识库
函数 更新知识库(knowledge_base, new_knowledge):
    // 将新检索到的知识合并到知识库中
    // 返回更新后的知识库

// 更新状态
```

```
函数 更新状态(state, step_answer, new_knowledge):  
    // 更新当前推理状态  
    // 返回更新后的状态  
  
// 满足结束条件  
函数 满足结束条件(state, knowledge):  
    // 判定是否已达到结论  
    // 返回布尔值  
  
// 生成最终答案  
函数 生成最终答案(state, knowledge):  
    // 根据当前状态和知识库生成答案  
    // 返回答案
```

## 3 Experiment

### 3.1 Self-consistency

1. input为： Q: When I was 6 my sister was half my age. Now I'm 70 how old is my sister? A:
  2. 将input分别输入GPT-4, CHATGPT, PALM三个模型中，输出为：
    - GPT-4: When you were 6 your sister was 3. Therefore, she is 3 years younger than you. So if you are 70 now, your sister is  $70 - 3 = 67$  years old.
    - CHATGPT : If when you were 6, your sister was half your age, that means when you were 6, your sister was 3 years old. The age difference between you and your sister is always the same, so when you're 70, your sister would be  $70 - 3 = 67$  years old.
    - PLAM: When you were 6 your sister was  $6 / 2 = 3$  years old. So your sister is  $70 - 3 = 67$  years old. The answer is 67.
  3. 最终结果为： The majority across the outputs indicates that the sister is 67 years old.
- 实验截图如下：



```
# Run the AI Config
from aiconfig.default_parsers.parameterized_model_parser import InferenceOptions

params = {}
inference_options = InferenceOptions()

response_completion = await config.run("majority-evaluator", params, inference_options, run_with_dependencies=True)
response = config.get_output_text("majority-evaluator")

print(response)

The majority across these outputs is the first one, with a calculation stating that the sister is 67 years old.
```

- 实验结果分析
  - 该实验是选择使用不同的LLM去生成多条推理链，该方法可以使得答案错误的概率较之只选用单一模型降低。
  - 然后该方法有着几个明显的弊端：
    1. 推理链的生成受限于大语言模型的数量；
    2. 多数即为答案的并不总是有效，有时候不同的语言模型会生成各不相同的答案，将没办法选出正确的答案是什么。

## 3.2 DIVERSE

- 数据集：GSM8K
- Input question: Pauline has 125 matchbox cars. They are all either convertibles, trucks, regular cars. 64% of them are regular cars. 8% are trucks. How many convertibles does she own?
- Truth: The trucks and regular cars make up 72% of her collection because  $64+8$  equals  $<<64+8>>72\%$  This means convertible make up 28% of her collection because  $100-72$  equals 28.%%She has 35 convertibles because 125 times .28 equals  $<<125*.28=35>>35\%$ #### 35
- sample0: There are 125 cars in total%%64% of them are regular cars so  $64/100*125=<<64/100*125=80>>80$  of them are regular cars%%8% of them are trucks so  $8/100*125=<<125-80-10=35>>35$  of them are convertibles%%#### 35
- 评价
  - 评价指标
    - random\_top1:是通过对数据集进行随机采样，然后计算随机选择的预测结果与真实答案匹配的频率来计算的。
    - voting\_top1\_accuracy:基于多数投票来评估预测性能的。对于每个样本，它计算预测结果中哪个答案被选为最多的次数，并将其与真实答案匹配。

- recall:是召回率，用于评估模型能够正确预测出真实答案的能力。

- 评价结果：

	训练集	测试集
random_top1	0.565	0.576
voting_top1_accuracy	0.783	0.782
recall	0.941	0.955

- 实验截图如下：

```
(microsoft.msra.dki.verifier_data_preparing) [ma-user src]$python verifier_data_prepare.py --generator_result_file executive/tra
in_dir --output_dir executive/train_preprocessed --split train --random_seed 233 --dataset_name GSM8K
/home/ma-user/anaconda3/envs/PyTorch-1.10.2/lib/python3.7/site-packages/requests/__init__.py:104: RequestsDependencyWarning: url
lib3 (1.26.12) or chardet (5.2.0)/charset_normalizer (2.0.12) doesn't match a supported version!
RequestsDependencyWarning)
100% | 20000/20000 [00:00<00:00, 200268.53it/s]
Total cases: 1000
Case 0's question: Pauline has 125 matchbox cars. They are all either convertibles, trucks, regular cars. 64% of them are regula
r cars. 8% are trucks. How many convertibles does she own?
Case 0's ground truth: The trucks and regular cars make up 72% of her collection because 64+8 equals <<64+8=72>>72%%This means c
onvertibles make up 28% of her collection because 100-72 equals 28.%%She has 35 convertibles because 125 times .28 equals <<125*
.28=35>>35%%### 35
Case 0's sample0: There are 125 cars in total%%64% of them are regular cars so 64/100*125=<<64/100*125=80>>80 of them are regula
r cars%%8% of them are trucks so 8/100*125=<<8/100*125=10>>10 of them are trucks%%The remaining cars are convertibles so 125-80-
10=<<125-80-10=35>>35 of them are convertibles%%### 35
***** Data statistics *****
random_top1: 0.565
voting_top1_accuracy: 0.783
recall: 0.941

100% | 1000/1000 [00:01<00:00, 592.95it/s]
100% | 1000/1000 [00:01<00:00, 976.66it/s]
100% | 1000/1000 [00:00<00:00, 192258.16it/s]
100% | 21000/21000 [00:00<00:00, 35624.84it/s]

(microsoft.msra.dki.verifier_data_preparing) root@autodl-container-bb8c4f8e12-349f008e:~# python DIVERSE/code/src/verifier_data_prepare.py --generator_result_file exec_dir/test_di
r --output_dir exec_dir/test_preprocessed --split dev --random_seed 233 --dataset_name GSM8K
100% | 99740/99740 [00:00<00:00, 233324.51it/s]
Total cases: 998
Case 0's question: Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at
the farmers' market daily for $2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?
Case 0's ground truth: Janet sells 16 - 3 - 4 = <<16-3-4=9>>9 duck eggs a day.%%She makes 9 * 2 = $<<9*2=18>>18 every day at the farmer's market.%%### 18
Case 0's sample0: She eats 3 and uses 4, so that is 7 eggs.%%16-7=<<16-7=9>>9 eggs to sell.%%9 eggs*$2=<<9*2=18>>$18/day%%### 18
***** Data statistics *****
random_top1: 0.5761523046092184
voting_top1_accuracy: 0.781563126232505
recall: 0.9549098196392786

100% | 998/998 [00:05<00:00, 187.83it/s]
100% | 998/998 [00:01<00:00, 505.74it/s]
100% | 998/998 [00:00<00:00, 88980.62it/s]
100% | 40918/40918 [00:01<00:00, 34566.20it/s]

(microsoft.msra.dki.verifier_data_preparing) root@autodl-container-bb8c4f8e12-349f008e:~#
```

- 实验结果分析

- 该方法效果很好，在多数投票来评估已经能显著提高准确率的情况下，再一次显著将准确率提高。在训练集和测试集均有着不错的表现。

### 3.3 IRCotT

- 实验步骤大致如下（详细步骤可参考./code/3\_IRCoT/README.md)

1. 配置实验环境：python 3.8.0
2. 准备数据集： hotpotqa , 2wikimultihopqa , musique , iirc
3. 准备prompt
4. 准备Retriever和LLM Servers

- SYSTEM: choose from ( ircot , ircot\_qa , oner , oner\_qa , nor\_qa )

- MODEL: choose from ( codex , flan-t5-xxl , flan-t5-xl , flan-t5-large , flan-t5-base , none )

#### 5. 输入问题

#### 6. 提取出问题中较为陌生的词语去检索

#### 7. 将检索出的词语带进问题，对问题进行推理

#### 8. 重复步骤6，7直到得出答案

### • 实验过程

What is the length of the river into which Pack Creek runs after it goes through the Spanish Valley?

A: ["Spanish Valley", "Spanish Creek (Plumas County, California)"]

A: Pack Creek runs through the Spanish Valley.

A: Exit? No.

A: ["Spanish Valley", "Spanish Creek (Plumas County, California)", "Polk Valley Run"]

A: Pack Creek runs through the Spanish Valley. After it goes through the Spanish Valley, it flows north northwestward into Moab Valley toward its confluence with the Colorado River.

A: Exit? No.

A: ["Spanish Valley", "Spanish Creek (Plumas County, California)", "Polk Valley Run", "Missouri River Valley"]

A: Pack Creek runs through the Spanish Valley. After it goes through the Spanish Valley, it flows north northwestward into Moab Valley toward its confluence with the Colorado River. The Colorado River is approximately 1,450 miles long.

A: Exit? No.

A: ["Spanish Valley", "Spanish Creek (Plumas County, California)", "Polk Valley Run", "Missouri River Valley", "Blue River (Colorado)", "Eagle River (Colorado)"]

A: Pack Creek runs through the Spanish Valley. After it goes through the Spanish Valley, it flows north northwestward into Moab Valley toward its confluence with the Colorado River. The Colorado River is approximately 1,450 miles long. So the answer is: 1,450 miles.

A:

```
["pid__c64c838ef73a481702ea6a0687937871__a8c2bbe4bd9ce63000deabd
0aee7f8ad",
"pid__2590840cb587f43dc1a25a7644854f01__11c071c82788906fa4ffa8977
b31473c",
"pid__7c024a4c8d5e5e889791e7eaaa8368d8__6fa3b2a7e67adb74b7edb26f6
d7d1f26",
"pid__f4ca3f6ccbe1ed2bbcc5376a4882f3c1__4272ea169eb66fad5baa9981ff
41abdc",
"pid__3497a744d9f0f980281031943126e461__a8f97b7902d3c7b16c06e75e
dd6f1f71",
"pid__d1dbb40ed2321e49981365ec25788dc3__d16b372ff9138d4892d1f9d8
8e0ccd74"]
```

Q: What is the length of the river into which Pack Creek runs after it goes through the Spanish Valley?

A: "Pack Creek runs through Spanish Valley and then into the Colorado River. The length of the Colorado River is approximately 1,450 mi. So the answer is: 1,450 mi."

A: 1,450 mi

Q: [EOQ]

S: 0.0

- 实验结果

```
*****
Experiment Name: ircot_qa_codex_hotpotqa
*****

python run.py summarize ircot_qa_codex_hotpotqa --instantiation_scheme
ircot_qa --prompt_set 1 --evaluation_path
processed_data/hotpotqa/dev_subsampled.jsonl

bm25_retrieval_count distractor_count      metric_value
0          2          "1" 60.5 | 64.6 | 61.4 | 100
1          2          "2" 60.9 | 64.0 | 62.8 | 100
2          2          "3" 60.1 | 63.7 | 61.2 | 100
3          4          "1" 55.6 | 58.8 | 57.3 | 100
4          4          "2" 57.6 | 61.1 | 58.7 | 100
```

5	4	"3"	55.3	58.8	56.7	100
6	6	"1"	55.6	59.1	56.8	100
7	6	"2"	52.2	55.6	53.5	100
8	6	"3"	48.9	52.0	50.2	100
9	8	"1"	53.5	56.4	54.7	100
10	8	"2"	52.7	55.6	53.6	100
...						

- 实验结果分析
  - 对于一些题目中蕴含复杂名词的问题，该方法可以借助检索来提高答案的生成。
  - 但是该方法的成功率不高，有相当一部分问题没办法检索出问题当中陌生名词的信息从而无法进行推理。该方法同样强烈依赖检索器的质量。

## 4 Reference

Self-Consistency Improves Chain of Thought Reasoning in Language Models

Making Large Language Models Better Reasoners with Step-Aware Verifier

Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions

- 
- 小组贡献占比：所有成员贡献均衡，均为20%