

# Lab4 CoT 实验

## 1 原理概述

### 1.1 思维链

Chain of Thought (COT) 推理是指一种自然语言处理 (NLP) 中的推理方式，其中模型生成一系列中间步骤或“思考链”，以解决复杂的问题或任务。在传统的 NLP 任务中，模型通常被要求直接给出答案，而在 COT 推理中，模型需要展示出它是如何一步一步推理出答案的。

COT 推理特别适用于解决需要多步逻辑推理或长时间记忆的任务，如解数学问题、回答多步推理问题等。这一方法类似于人类在解决问题时先内部推理出问题的答案，再把最终答案说出来。Chain of Thought 推理通过这种方式提高了解决复杂任务的准确性，并使预训练语言模型的输出更加透明和可理解。

参考论文：Large Language Models are Zero-Shot Reasoners

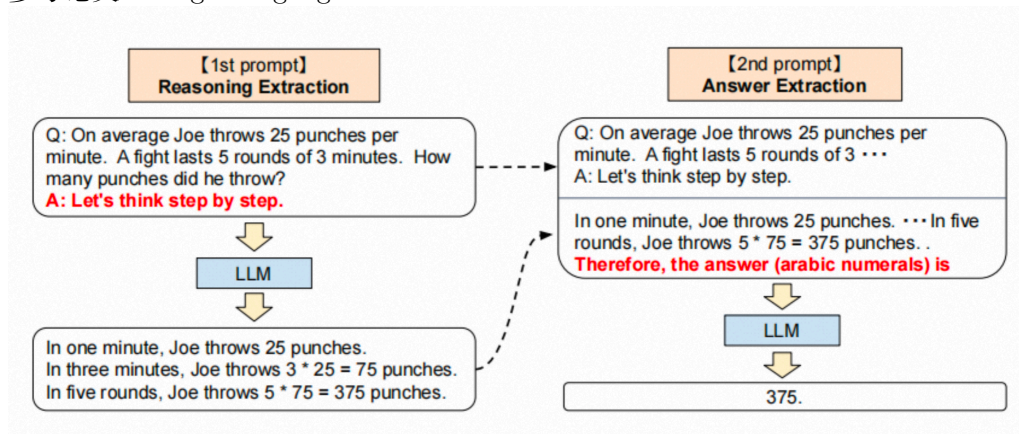


图 1: CoT 过程

### 1.2 CoT 例子

#### 1.2.1 算术问题解答过程说明

- **第一步提示 ( $X_0$ ):** 采用模板“Q: [X]. A: [T]”，其中 [X] 是问题槽位，[T] 是引导推理过程的触发句，默认为“Let's think step by step”
- **后续生成句子 Z:** 将第一步提示  $X_0$  发送到 LLaMA 中，以产生后续句子  $z$ 。
- **第二步提示:** 结合  $X_0$  和  $Z$ ，构成新的提示“[ $X_0$ ] [ $Z$ ] [A]”，其中 [A] 是触发模型输出答案的模板。
- **最终结果:** 将完成的提示作为输入送入 LLaMA，得到答案预测 sentence  $\hat{y}$ 。
- **准确率计算:** 对算术问题，提取预测 sentence  $\hat{y}$  中的数字作为预测答案，以计算模型准确率。

#### 1.2.2 例子：解决一个简单的数学问题

- **问题:** 一个园丁有 24 朵花，他平均分给了 4 个孩子，每个孩子得到了多少朵花？

- **第一步提示 ( $X_0$ ):** 采用模板 “Q: 一个园丁有 24 朵花, 他平均分给了 4 个孩子, 每个孩子得到了多少朵花? A: Let’s think step by step.”
- **后续生成句子 ( $Z$ ):** 发送  $X_0$  到 LLaMA 模型, 得到  $Z$ : “To find out how many flowers each child gets, we need to divide the total number of flowers by the number of children.”
- **第二步提示:** 结合  $X_0$  和  $Z$ , 构成新的提示: “Q: 一个园丁有 24 朵花, 他平均分给了 4 个孩子, 每个孩子得到了多少朵花? A: Let’s think step by step. To find out how many flowers each child gets, we need to divide the total number of flowers by the number of children. [A]” 这里的 [A] 代表待插入答案的模板部分。
- **最终结果:** 完成的提示作为输入发送到 LLaMA 模型。得到预测句子 ( $\hat{y}$ ): “Each child gets 24 divided by 4, which is 6 flowers.”
- **准确率计算:** 从  $\hat{y}$  中提取的数字是 6。这是正确答案 ( $24 \div 4 = 6$ ), 表示模型解决了问题。对于模型准确率的评估, 这将记为正确解答。通过多项此类问题的正确解答比例, 我们得到模型的准确率。

### 1.2.3 常识问答部分过程说明

- **第一步提示 ( $X_0$ ):** 同上, 采用模板 “Q: [X]. A: [T]”。
- **后续生成句子  $Z$ :** 同上, 将第一步提示发送到 LLaMA 中, 以产生后续句子  $z$ 。
- **第二步提示:** 结合  $X_0$  和  $Z$ , 构成新的提示 “[X0] [Z] [A]”, 但 [A] 的模板稍有不同, 更改为 “Therefore, among A through E, the answer is”。
- **最终结果:** 同上, 将完成的提示作为输入送入 LLaMA, 得到答案预测 sentence  $\hat{y}$ 。
- **准确率计算:** 对常识问答, 提取预测 sentence  $\hat{y}$  中遇到的第一个大写字母作为预测答案, 以计算模型准确率。

### 1.2.4 例子: 解决一个常识选择题

- 假设问题是: “哪位科学家首次提出了相对论理论? 选项如下: A) 尼古拉·特斯拉 B) 艾萨克·牛顿 C) 阿尔伯特·爱因斯坦 D) 尼尔斯·玻尔”
- **第一步提示 ( $X_0$ ):** 采用模板 “Q: 哪位科学家首次提出了相对论理论? 选项如下: A) 尼古拉·特斯拉 B) 艾萨克·牛顿 C) 阿尔伯特·爱因斯坦 D) 尼尔斯·玻尔 A: Let’s think step by step.”
- **后续生成句子  $Z$ :** 发送  $X_0$  到 LLaMA 模型, 得到  $Z$ : “The theory of relativity, which revolutionized physics, was introduced by a scientist known for his work in theoretical physics during the early 20th century. This theory includes concepts like time dilation and the equivalence of mass and energy.”
- **第二步提示:** 结合  $X_0$  和  $Z$ , 构成新的提示: “Q: 哪位科学家首次提出了相对论理论? 选项如下: A) 尼古拉·特斯拉 B) 艾萨克·牛顿 C) 阿尔伯特·爱因斯坦 D) 尼尔斯·玻尔 A: Let’s think step by step. The theory of relativity, which revolutionized physics, was introduced by a scientist known for his work in theoretical physics during the early 20th century. This theory includes concepts like time dilation and the equivalence of mass and energy. Therefore, among A through D, the answer is”

- **最终结果：**完成的提示作为输入发送到 LLaMA 模型。模型可能会返回预测句子  $\hat{y}$ : “Therefore, among A through D, the answer is C, which stands for 'Albert Einstein'.”
- **准确率计算：**从预测句子  $\hat{y}$  中提取出的第一个大写字母是 C。由于 C 代表的是阿尔伯特·爱因斯坦，这是正确的答案，因为爱因斯坦确实首次提出了相对论理论。计算模型准确率时，这将被记为一个正确答案。如果我们有一系列类似的问题和答案，可以通过正确答案的比例来计算模型的准确率。

## 2 实验内容

### 1. 数据预处理：

- **GSM8K 数据集：**
  - a) 可在 `cot/data` 中找到，格式为 `parquet`，只处理前 50 条记录。
  - b) 输入集：仅使用列 `question` 作为输入。
  - c) Label 集：提取 `answer` 中 “#####” 之后的部分作为正确答案。
- **CommonsenseQA 数据集：**
  - a) 可在 `cot/data` 中找到，处理前 50 条记录。
  - b) 输入集：将 `question`、`choices` 和 `text` 进行拼接。
  - c) Label 集：使用 `answerKey` 作为正确答案。

### 2. COT 推理过程：

- 参考 `llam_infer` 代码
  - a) 代码可在以下链接找到：[https://github.com/mindspore-courses/step\\_into\\_llm/blob/master/Season2.step\\_into\\_llm/04.LLaMA/llama\\_infer.py](https://github.com/mindspore-courses/step_into_llm/blob/master/Season2.step_into_llm/04.LLaMA/llama_infer.py) 此代码提供了如何调用 LLaMA 接口进行推理
- **算术数据部分：**
  - a) 第一步提示  $X_0$ ：采用模板 “Q: [X]. A: [T]”，其中 [X] 是问题的槽位，[T] 用于引导推理过程的触发句（默认为 “Let’s think step by step”）。
  - b) 后续生成句子  $Z$ ：将第一步提示发送到 LLaMA 中，产生后续句子  $z$ 。
  - c) 第二步提示：结合第一步提示  $X_0$  与  $Z$ ，构成新的提示 “[X0] [Z] [A]”，[A] 为触发模型输出答案的模板（默认为 “Therefore, the answer (arabic numerals) is”）。
  - d) 最终结果：将完成的提示作为输入送入 LLaMA，获得答案预测 sentence  $\hat{y}$ 。
  - e) 准确率计算：对算术问题，提取预测 sentence  $\hat{y}$  中的数字作为预测答案，以计算模型准确率。
- **常识问答部分：**
  - a) 第一步提示  $X_0$ ：采用模板 “Q: [X]. A: [T]”，其中 [X] 是问题的槽位，[T] 用于引导推理过程的触发句（默认为 “Let’s think step by step”）。
  - b) 后续生成句子  $Z$ ：将第一步提示发送到 LLaMA 中，产生后续句子  $z$ 。
  - c) 第二步提示：结合第一步提示  $X_0$  与  $Z$ ，构成新的提示 “[X0] [Z] [A]”，[A] 为触发模型输出答案的模板，此时触发句 [A] 格式稍有不同，应修改为 “Therefore, among A through E, the answer is”。
  - d) 最终结果：将完成的提示作为输入送入 LLaMA，获得答案预测 sentence  $\hat{y}$ 。

e) 准确率计算：遇到的第一个大写字母将作为模型的预测答案，以此来计算模型准确率。

3. **实验探究：**尝试不同的模板构建方法以及不同的触发句，通过准确率选择出最佳的 COT 生成方式。