✏️

# 6. (2022)AUTOMATIC CHAIN OF THOUGHT PROMPTING IN LARGE LANGUAGE MODELS

| ⊙ Created by | 🖼️ 哲睿 張 |
| --- | --- |
| 🕐 Created time | @February 22, 2024 6:45 PM |
| ☰ Tags | |
| 🔗 https://www.youtube.com/watch?v=l1woW9_vR1c | https://zhuanlan.zhihu.com/p/618904090 |
| 🔗 https://www.youtube.com/watch?v=l1woW9_vR1c (1) | https://cloud.tencent.com/developer/article/2321413 |

> https://prod-files-secure.s3.us-west-2.amazonaws.com/0b0c1a86-b713-4c99-9c35-4c26e958a80d/b420a3d9-e32f-4c22-98c5-56c8fde6b95f/AutoCoT.pdf

## 1. Abstract

**利用GPT-3進行自動化生成CoT，效果希望超越手動CoT。**
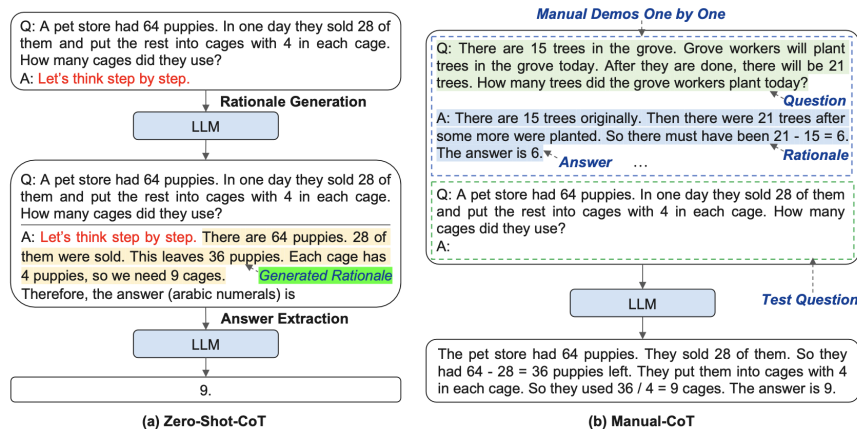
> 💡 Let's not just think step by step, but also one by one.

**發現多樣性的CoT可以減少錯誤CoT帶來的效果影響。**

**Auto-CoT步驟**

> 💡 1. 將給定數據集的問題劃分為幾個群組。
> 2. 從每個群組中選擇一個代表性問題,並使用簡單的啟發式方法使用零射 CoT 生成其推理鏈。

# 2. Related work



Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?
A: Let's think step by step.

**Rationale Generation**

LLM

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?
A: Let's think step by step. There are 64 puppies. 28 of them were sold. This leaves 36 puppies. Each cage has 4 puppies, so we need 9 cages. *Generated Rationale* Therefore, the answer (arabic numerals) is

**Answer Extraction**

LLM

9.

**(a) Zero-Shot-CoT**

*Manual Demos One by One*

Q: There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?
*Question*
A: There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been 21 - 15 = 6. The answer is 6. *Answer* … *Rationale*

Q: A pet store had 64 puppies. In one day they sold 28 of them and put the rest into cages with 4 in each cage. How many cages did they use?
A:

LLM

*Test Question*

The pet store had 64 puppies. They sold 28 of them. So they had 64 - 28 = 36 puppies left. They put them into cages with 4 in each cage. So they used 36 / 4 = 9 cages. The answer is 9.

**(b) Manual-CoT**

(兩種範式: Zero-shot-CoT ; Manual-CoT)

作者回顧2種相關任務:

1. CoT prompting (又分成zeroshot和manial CoT)

2. In context learning (ICL)

   a. 有學者質疑(2022) ICL的正確與否並不會影響表現太多,是因為任務往往是標準分類 <Input, Output>映射

   b. 對於較複雜的映射(<input, rationale, output>)出現錯誤時,效能急遽下降

# 3. Challenge of AutoCoT

不同標記者(annotator)寫出不同的demo,在符號推理任務上差距甚至達到28% → 怎麼寫範例很重要。

## 3.1 挑出MultiArith資料集來了解為什麼Retrieval-Q-CoT比Random-Q-CoT差

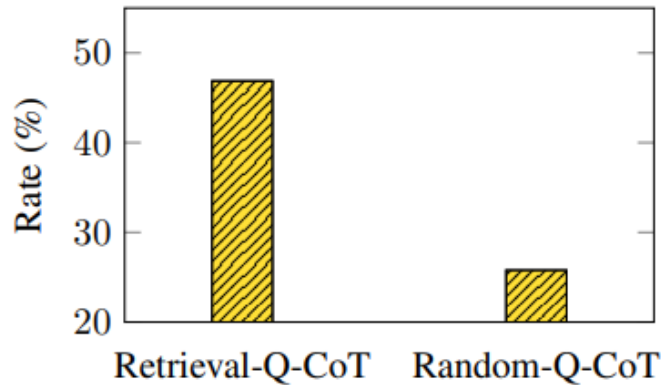| Method | MultiArith | GSM8K | AQuA |
|---|---|---|---|
| Zero-Shot-CoT | 78.7 | 40.7 | 33.5 |
| Manual-CoT | **91.7** | 46.9 | 35.8† |
| Random-Q-CoT | 86.2 | 47.6† | 36.2† |
| Retrieval-Q-CoT | 82.8 | **48.0†** | **39.7†** |



Figure 2: Unresolving Rate.

💡 zero-shot-CoT中(只有"Let;s think strep by step"這個咒語)，錯誤率21.3％
(128/600)，而當額外加入Retrieval-Q-CoT和Random-Q-CoT後，

依然失敗的占比在Figure.2。Retrieval-Q-CoT比Random-Q-CoT更差的原因是因為
Retrieval-Q-CoT使用相似度方法，造成"一步錯步步錯"的概念。

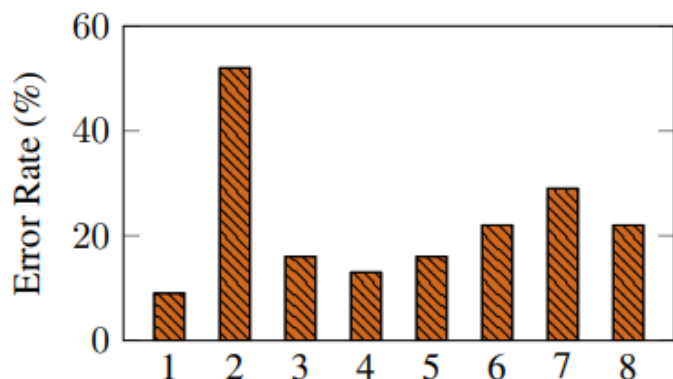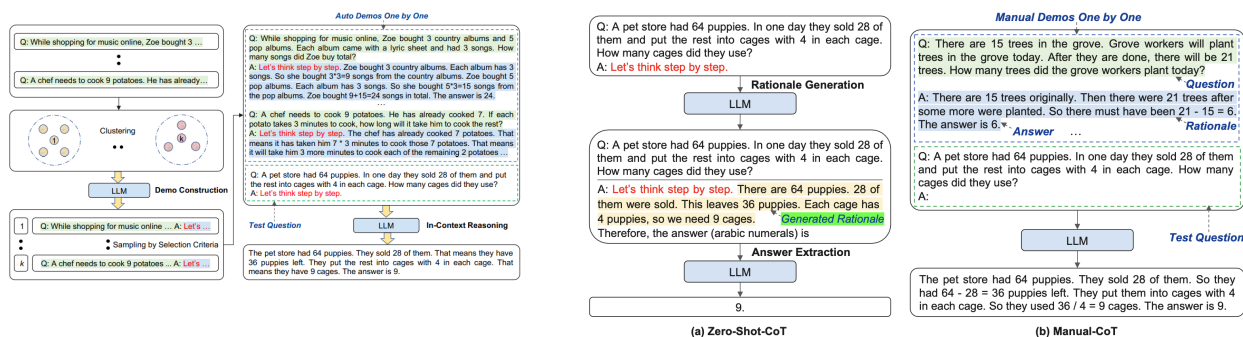| Retrieval-Q-CoT | Random-Q-CoT |
|---|---|
| **Q:** A chef needs to cook 9 potatoes. He has already cooked 7. If each potato takes 3 minutes to cook, how long will it take him to cook **the rest**? <br> **A:** Let's think step by step. The chef has already cooked 7 potatoes. That means it has taken him 7 * 3 minutes to cook those 7 potatoes. That means it will take him 3 more minutes to cook each of the remaining 2 potatoes. That means it will take him a total of 7 * 3 + 2 * 3 = 21 minutes to cook all 9 potatoes. The answer is 21. ✗ <br> **Q:** A chef needs to cook 13 potatoes. He has already cooked 5. If each potato takes 6 minutes to cook, how long will it take him to cook **the rest**? <br> **A:** Let's think step by step. The chef has already cooked 5 potatoes. That means it has taken him 30 minutes so far. There are 13 potatoes in total. 5 have been cooked so far. That means there are 8 left to cook. Each potato takes 6 minutes to cook. That means it will take the chef 48 minutes to cook the 8 remaining potatoes. In total, it will take the chef 78 minutes to cook all 13 potatoes. The answer is 78. ✗ | **Q:** Wendy uploaded 79 pictures to Facebook. She put 44 pics into one album and put the rest into 5 different albums. How many pictures were in each album? <br> **A:** Let's think step by step. Wendy uploaded 79 pictures to Facebook. She put 44 pictures into one album. That means she put the rest of the pictures into 5 different albums. That means each album would have (79-44)/5 = 15 pictures in each album. The answer is 15. ✓ <br> **Q:** At a restaurant each adult meal costs $5 and kids eat free. If a group of 15 people came in and 8 were kids, how much would it cost for the group to eat? <br> **A:** Let's think step by step. There are 8 kids, so that means there are 7 adults. Each adult meal costs $5, so that means that the 7 adults will cost $5*7=$35. Now we just need to add the two numbers together. $35+8=43$. So it would cost $43 for the group to eat. The answer is 43. ✗ |
| **Q:** A chef needs to cook 15 potatoes. He has already cooked 8. If each potato takes 9 minutes to cook, how long will it take him to cook **the rest**? <br> **A:** Let's think step by step. The chef has already cooked 8 potatoes. That means it has taken him 8 * 9 minutes to cook those 8 potatoes. That means it will take him 9 minutes to cook each of the remaining 7 potatoes. That means it will take him a total of 8 * 9 + 7 * 9 = 135 minutes to cook all 15 potatoes. The answer is 135. ✗ | **Q:** A chef needs to cook 15 potatoes. He has already cooked 8. If each potato takes 9 minutes to cook, how long will it take him to cook **the rest**? <br> **A:** Let's think step by step. The chef needs to cook 15 potatoes. He has already cooked 8. That means he has to cook 15-8=7 more potatoes. Each potato takes 9 minutes to cook. That means it will take him 9*7=63 minutes to cook **the rest** of the potatoes. The answer is 63. ✓ |

## 3.2證明Retrieval-Q-CoT失敗的原因是因為有"頻繁錯誤群集"

Figure 3: Clusters of similar questions.

💡 Cluster.2可能是因為缺乏問題解決的常用面相，所以才會解決不了問題

# 4. AutoCoT map(示意圖)



## 4.1 Question Clustering問題分群

- 先用Sentence-BERT計算每個question的向量長度，並且統一所有問題的vector大小。

- K-means將questions分群

## 4.2 Demonstration Sampling

- 假設現在Question set中有k個cluster，產生出[Q: q(i)j. A: [P]]

- 丟入LLM產生出rationale，形成[Q: q(i)j , A: r(i)j。a(i)j] →一個完整demonstration ( 限制:60tokens內,rationale在5個步驟以內)

**Algorithm 1** Cluster

**Require:** A set of questions $\mathcal{Q}$ and the number of demonstrations $k$
**Ensure:** Sorted questions $\mathbf{q}^{(i)} = [q_1^{(i)}, q_2^{(i)}, \ldots]$ for each cluster $i$ $(i = 1, \ldots, k)$
1: **procedure** CLUSTER($\mathcal{Q}, k$)
2:     **for** each question $q$ in $\mathcal{Q}$ **do**
3:         Encode $q$ by Sentence-BERT
4:     Cluster all the encoded question representations into $k$ clusters
5:     **for** each cluster $i = 1, \ldots, k$ **do**
6:         Sort questions $\mathbf{q}^{(i)} = [q_1^{(i)}, q_2^{(i)}, \ldots]$ in the ascending order of the distance to the cluster center
7:     **return** $\mathbf{q}^{(i)}$ $(i = 1, \ldots, k)$

**Algorithm 2** Construct

**Require:** Sorted questions $\mathbf{q}^{(i)} = [q_1^{(i)}, q_2^{(i)}, \ldots]$ for each cluster $i$ $(i = 1, \ldots, k)$, empty demonstration list $\mathbf{d}$
**Ensure:** Demonstration list $\mathbf{d} = [d^{(1)}, \ldots, d^{(k)}]$
1: **procedure** CONSTRUCT($\mathbf{q}^{(i)}, \ldots, \mathbf{q}^{(k)}$)
2:     **for** each cluster $i = 1, \ldots, k$ **do**
3:         **for** each question $q_j^{(i)}$ in $\mathbf{q}^{(i)}$ **do**
4:             Generate rationale $r_j^{(i)}$ and answer $a_j^{(i)}$ for $q_j^{(i)}$ using Zero-Shot-CoT
5:             **if** $q_j^{(i)}, r_j^{(i)}$ satisfy selection criteria **then**
6:                 Add $d^{(i)} = [\text{Q: } q_j^{(i)}, \text{A: } r_j^{(i)} \circ a_j^{(i)}]$ to $\mathbf{d}$
7:                 **break**
8:     **return** $\mathbf{d}$
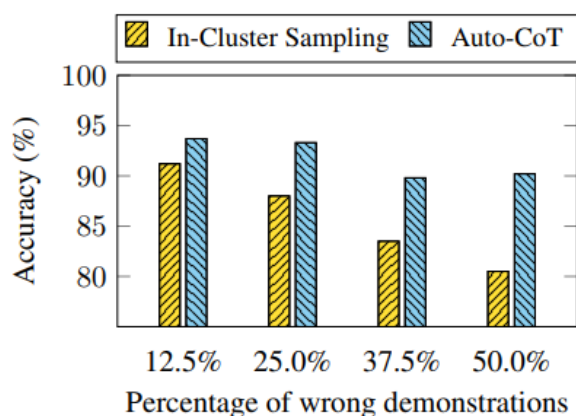
# 5. Experiments

**任務類型:**

(i) 算術推理 　(ii)符號推理 　(iii)常識推理

模型: GPT-3 (text-davinci-002 version 175B)

| Model | Arithmetic | | | | | | Commonsense | | Symbolic | |
|---|---|---|---|---|---|---|---|---|---|---|
| | MultiArith | GSM8K | AddSub | AQuA | SingleEq | SVAMP | CSQA | Strategy | Letter | Coin |
| Zero-Shot | 22.7 | 12.5 | 77.0 | 22.4 | 78.7 | 58.8 | 72.6 | 54.3 | 0.2 | 53.8 |
| Zero-Shot-CoT | 78.7 | 40.7 | 74.7 | 33.5 | 78.7 | 63.7 | 64.6 | 54.8 | 57.6 | 91.4 |
| Few-Shot | 33.8 | 15.6 | 83.3 | 24.8 | 82.7 | 65.7 | **79.5** | **65.9** | 0.2 | 57.2 |
| Manual-CoT | 91.7 | 46.9 | 81.3 | 35.8 | 86.6 | 68.9 | 73.5 | 65.4 | 59.0 | 97.2 |
| Auto-CoT | **92.0** | **47.9** | **84.8** | **36.5** | **87.0** | **69.5** | 74.4 | 65.4 | **59.7** | **99.9** |

Auto-CoT全面性優於manual-CoT，原因是手動的成本很高，設計者在設計dataset的CoT時不會一個個demo(在算數dataset中5/6的demo都是同一個)。

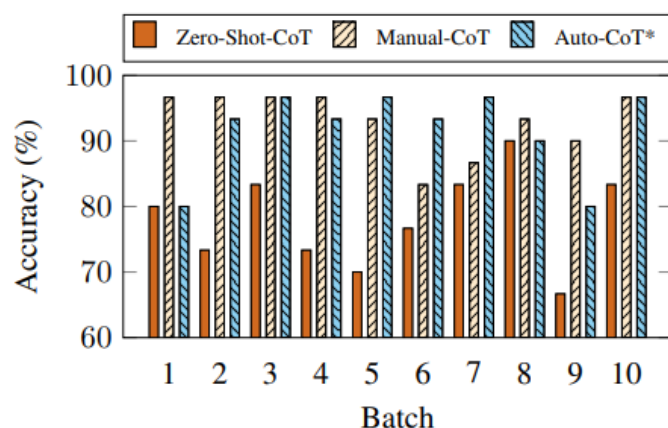比較起來，Auto-CoT比manual-CoT更具彈性和任務的適應性。

## 錯誤的Demo造成的影響

(In-Cluster Sampling: 在同一集群中隨機抽樣問題)

## 更具挑戰的任務

說明: 當Dtaset不是一次完整丟進來，而是批次陸續丟入。



# 6. Conclusion

LLMs已經展示了在CoT提示下的推理能力。Manual-CoT的卓越性能取決於示範的手工設計。為了消除這種手工設計，我們提出了Auto-CoT來自動構建示範。

它通過多樣性抽樣問題並生成推理鏈來構建示範。對於十個公共基準推理數據集的實驗結果表明，使用GPT-3，Auto-CoT始終與需要手工設計示範的CoT範式的性能相匹配或超越。