# Inference with CoT Based on LLM (Topic 2)

## 1. Project Introduction

### 1.1 Development environment and requirements for system

- **Development environment:** Google Colab / Ubuntu 22.4

- **Framework:** pytorch 2.3.0 (cu121)

- **Third-party Package:**

  - transformers 4.41.2

  - modelscope 1.15.0

### 1.2 Main step of the project

**a.GOAL**:Improve inference with CoT accuracy of LLM through different methods

**b.Dataset**: AQuA:https://github.com/deepmind/AQuA

**c.Method**: We have implemented three methods to improve the inference with CoT accuracy of LLM:

- Multi-vote

- LoRA fine-tuning

- Multi-path step-aware inference

**d.Evaluation**: We evaluate the performance of the model by the accuracy of the inference with CoT.

## 2. Technical Details

### 2.1 Chain of Thought (CoT)

The Chain-of-Thought (CoT) methodology represents the cognitive reasoning process during language understanding. In LLMs, CoT utilizes logical inference and concept associations to comprehend and generate text. By constructing a chain of thought, LLMs can more accurately understand and answer questions, as they are able to connect various components of the question into a coherent reasoning process.
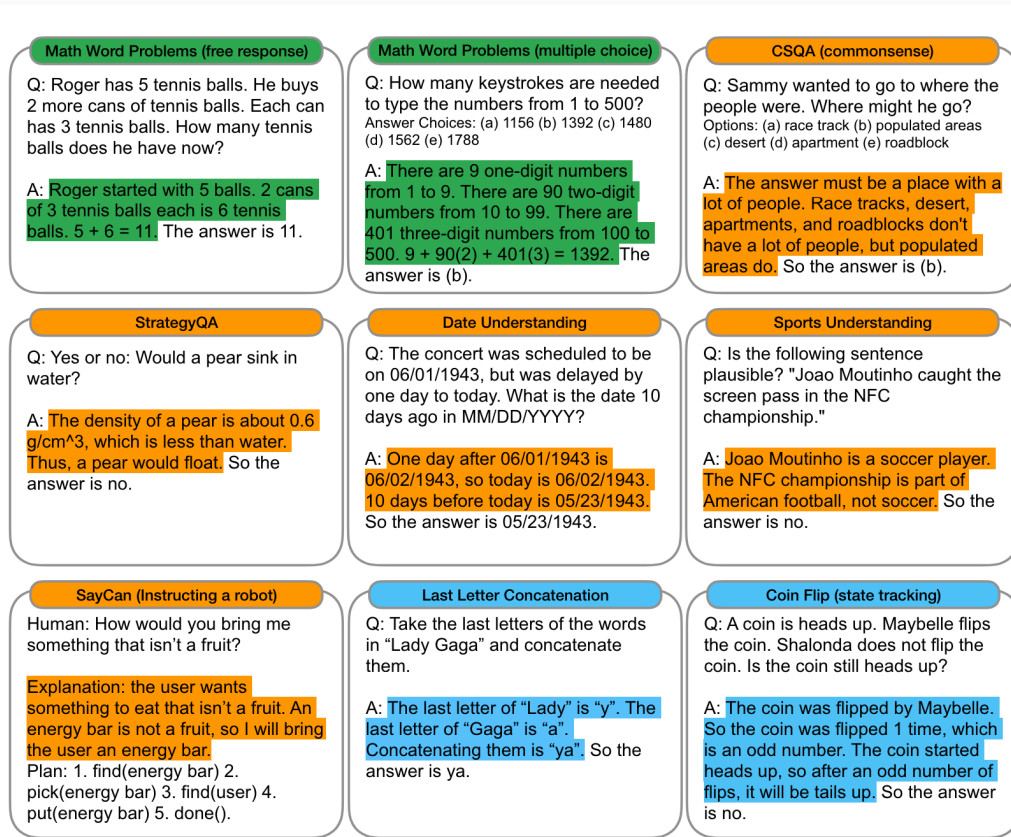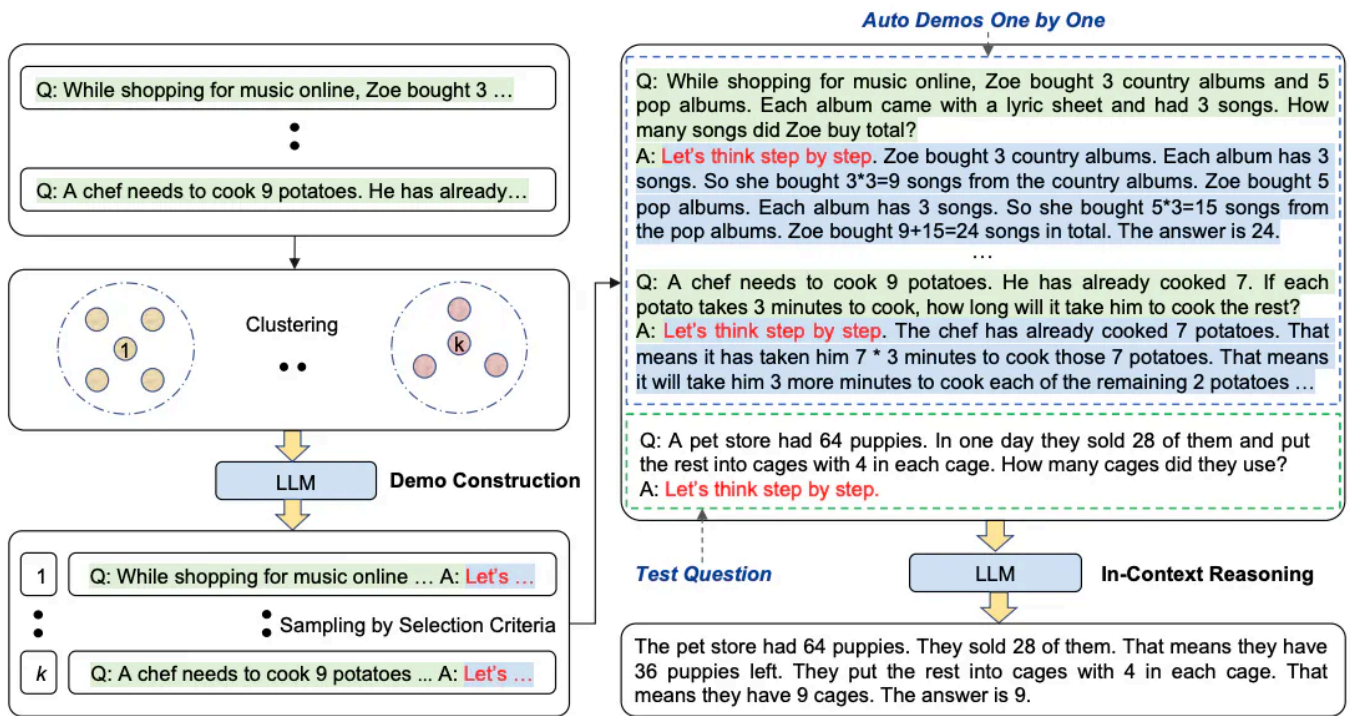


Figure 3: Examples of ⟨input, chain of thought, output⟩ triples for arithmetic, commonsense, and symbolic reasoning benchmarks. Chains of thought are highlighted. Full prompts in Appendix G.

COT can be classified into two types: Few-shot CoT based on manually annotated examples, and Zero-shot CoT without manually annotated examples. Few-shot CoT uses human-annotated examples to construct chains of thought, while Zero-shot CoT attempts to generate chains of thought autonomously. The introduction of chains of thought enables large-scale language models to excel in handling complex language tasks such as arithmetic reasoning, common-sense reasoning, and symbolic reasoning.

## Reasoning Process:

1. **Initial Prompt ($X_0$):** This step involves formulating an initial prompt using a predefined template structure. By separating the question and answer segments, the template ensures clarity and guides the subsequent reasoning process. The trigger sentence within the prompt sets the tone for how the reasoning will unfold, emphasizing a step-by-step approach.

2. **Generating Subsequent Sentences ($Z$):** After presenting the initial prompt to the LLM, subsequent sentences $Z$ are generated. These sentences provide additional context, reasoning steps, or elaborations based on the question posed. By generating multiple sentences, the LLM has the opportunity to explore various lines of thought and provide a comprehensive rationale.

3. **Subsequent Prompt:** The subsequent sentences $Z$ are integrated with the initial prompt $X_0$ to create a new prompt. This combined prompt serves as the basis for further interaction with the LLM, guiding it towards generating a coherent answer. By incorporating the additional context provided by $Z$, the subsequent prompt enriches the overall reasoning process.

4. **Final Result:** With the completed prompt in hand, the LLM generates the predicted answer sentence $\hat{y}$. This sentence represents the culmination of the reasoning process, encapsulating the LLM's understanding and interpretation of the given question and context.

5. **Accuracy Calculation:** Depending on the nature of the task (e.g., arithmetic problems or common-sense question answering), the accuracy of the predicted answer is calculated. This step involves extracting relevant information from the predicted sentence $\hat{y}$ and comparing it to the ground truth to evaluate the model's performance.



## 2.2 Multi-vote

- Specific implementation can be found in code/multi-vote.ipynb

### 2.2.1 Single-vote

The single-vote method is a straightforward approach to inference with CoT. It involves generating a single reasoning path using a predefined prompt and evaluating the model's answer based on this path.

- The main steps of the single-vote method are as follows:

    1. **Initial Prompt:** Construct an initial prompt using a predefined template structure.

    2. **Inference:** Generate text based on the input question and initial prompt using the LLM.

    3. **Decode and Process:** Decode the model outputs, extract the generated text, and process it to obtain the final answer.

    4. **Evaluation:** Compare the predicted answer with the ground truth to evaluate the model's performance.

**Algorithm 1** Inference and Prediction Process

```
1: procedure INFERENCE(model, tokenizer, question, max_new_tokens)
2:     input_ids ← tokenizer.encode(question, return_tensors="pt").to(model.device)
3:     outputs ← model.generate(input_ids, max_new_tokens)
4:     return decode(outputs, tokenizer, len(input_ids))
5: end procedure
6: predictions ← []
7: prompts ← ["Let's think step by step.", "Therefore, among A through E,
   the label of the answer is"]
8: for questions, answers in data_loader do
9:     for question, answer in zip(questions, answers) do
10:        input_1 ← "Q: " + question + ". " + prompts[0]
11:        output_1 ← INFERENCE(model, tokenizer, input_1, 256)
12:        input_2 ← input_1 + output_1 + prompts[1]
13:        output_2 ← INFERENCE(model, tokenizer, input_2, 10).replace(input_2, "")
14:        prediction ← FINDALL('[A-E]', output_2)[-1] or None
15:        Append prediction to predictions
16:    end for
17: end for
```

## 2.2.2 Multi-vote

The multi-vote method extends the single-vote approach by generating multiple reasoning paths using different prompts or hyperparameters. By leveraging multiple paths, the model can explore diverse reasoning strategies and improve the robustness of its predictions.

The main steps of the multi-vote method are as follows:

1. **Multiple Prompts:** Define multiple prompts or hyperparameters to generate different reasoning paths.

2. **Inference:** Generate text based on the input question and each prompt using the LLM.

3. **Voting:** Aggregate the predictions from multiple paths and select the final answer based on majority voting.

4. **Evaluation:** Compare the predicted answer with the ground truth to evaluate the model's performance.

**Algorithm 2** Multi Infer and Vote

```
1: procedure MULTIINFERANDVOTE(model, tokenizer, question, first_prompts, second_prompts)
2:     assert len(first_prompts) == len(second_prompts)
3:     votes ← []
4:     for i ← 0 to len(first_prompts) - 1 do
5:         first_input ← "Q : " + question + first_prompts[i]
6:         first_output ← INFERENCE(model, tokenizer, first_input, 350)
7:         second_input ← first_output + second_prompts[i]
8:         second_output ← INFERENCE(model, tokenizer, second_input, 10).replace(second_input,
   "")
9:         ans ← FINDALL('[A-E]', second_output)[-1]
10:        if ans then
11:            Append ans to votes
12:        end if
13:    end for
14:    if votes is empty then
15:        return None
16:    else
17:        return max(set(votes), key=votes.count)
18:    end if
19: end procedure
```
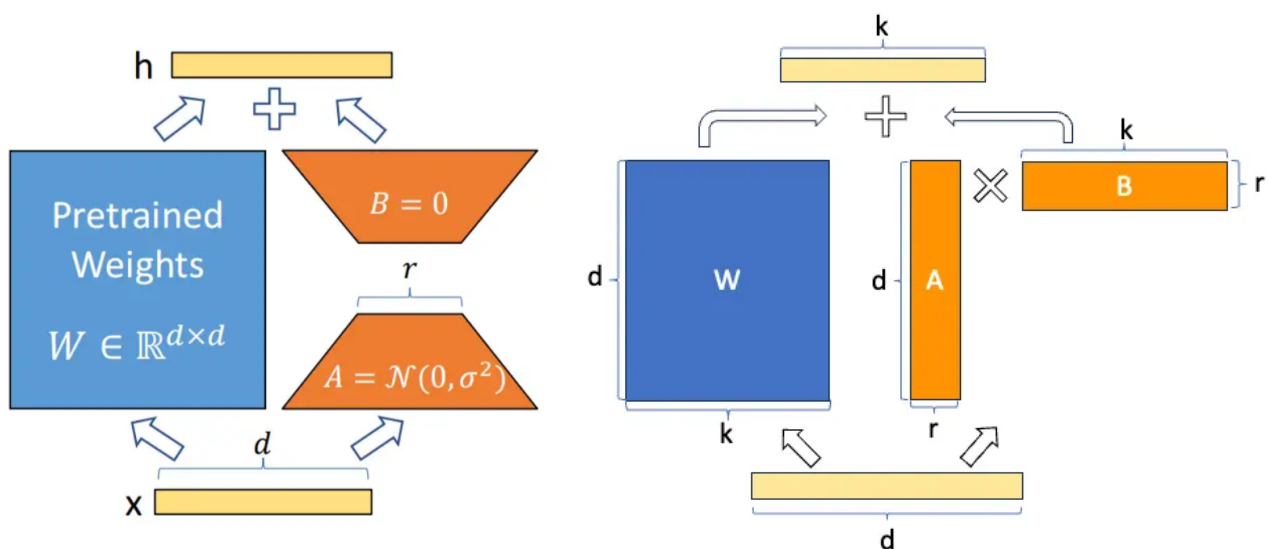
## 2.3 LoRA fine-tuning

- Specific code can be found in `code/lora_tune.ipynb`

### 2.3.1 Introduction for LoRA

Ordinarily, fine-tuning involves training the model on the new dataset, updating weights on the basis of pre-trained model directly.

Most, if not all, of the model's parameters are updated during training. It requires substantial computational resources (memory and processing power) and time, especially for large models. It may also require a large amount of data to avoid overfitting and to fine-tune the model effectively. High performance can often be achieved on the target task due to the comprehensive update of model parameters.



Low-Rank Adaptation (LoRA) is a parameter-efficient fine-tuning technique. Its core idea is to perform an implicit low-rank transformation on the weight matrices of large models, introducing additional trainable parameters. Here's how LoRA works:

1. **Model Loading**: Load a pre-trained model.

2. **LoRA Layers Addition**: Insert low-rank matrices into each layer of the model. These matrices have significantly fewer parameters compared to the original model layers.

3. **Training**: Train the model with the low-rank matrices, keeping the original model parameters frozen. **Only the LoRA parameters are updated.**

### 2.3.2 Characteristics of LoRA:

- **Parameter Update**: Only a small fraction of the model's parameters (the low-rank matrices) are updated during training.

- **Resource Efficiency**: Much less computationally intensive as fewer parameters are being optimized. Requires less memory and processing power.

- **Data Efficiency**: Often requires less data compared to ordinary fine-tuning since fewer parameters are being tuned.

- **Performance**: Can achieve competitive performance with ordinary fine-tuning, especially in resource-constrained environments or when the available data is limited.

Some features comparing are as below:

| Feature | Ordinary Fine-Tuning | LoRA Fine-Tuning |
| --- | --- | --- |
| Parameter Update | Most or all model parameters | Small fraction of parameters (low-rank matrices) |
| Resource Requirement | High (memory, computation) | Low (memory, computation) |
| Data Requirement | Often requires large datasets | Can work well with smaller datasets |
| Training Time | Longer | Shorter |
| Complexity | Simpler but resource-heavy | Slightly more complex but efficient |
| Performance | Typically high | Competitive, especially in constrained settings |

LoRA fine-tuning offers a more efficient alternative by updating only a small set of parameters, making it ideal for scenarios where computational resources and data are limited.

### 2.3.3 Our implementation

Hyperparameters of inference are as below (the values of hyperparameters are not specific here):

```
1  hyperparameters = {
2      "do_sample": True,
3      "max_new_tokens": 100,
4      "temperature": 0.7,
5      "top_k": 80,
6      "top_p": 0.9,
7      "num_return_sequences": 1,
8      "pad_token_id": pad_token_id
9  }
```

The prototype of inference function `inference` is as below:

```
1  def inference(model, tokenizer, question, max_new_tokens, num_beams=5, temperature=0.7, top_p=0.9)
```

Here are the meanings of some hyperparameters and arguments of inference.

| Parameter Name | Description | Purpose | Recommendation |
|---|---|---|---|
| `temperature` | Temperature | Controls the randomness of the generated text. Higher temperature (e.g., > 1) results in more random and diverse output, while lower temperature (e.g., <= 0.5) makes the model generate more deterministic and common text. | Use a higher temperature for more creativity and diversity. Use a lower temperature for accuracy and consistency. |
| `top_k` | Top-K Sampling | Considers only the top K most probable words when predicting the next word. | Smaller K values (like 50 or 100) reduce randomness and improve generation quality but may decrease diversity. Larger K values increase diversity but may lower quality. Adjust based on experimentation. |
| `top_p` | Top-P Sampling | Also known as nucleus sampling, restricts the next word selection to the set of words whose cumulative probability reaches P. | Typically set between 0.9 and 1 to balance quality and diversity. A value closer to 1 results in outputs similar to greedy search (always choosing the highest probability word), while a smaller value increases diversity. |
| `num_beams` | Num Beams | In Beam Search, determines the number of best candidate sequences explored simultaneously. | Set a larger value (e.g., 20) for higher output quality at the cost of increased computational resources. For limited resources or speed requirements, use a smaller value. |
| `max_new_tokens` | Max Tokens | Limits the maximum length of the generated text. | Set according to the application to ensure the content is not too long or too short. For example, use a smaller one for summaries and larger one for storytelling. |
| `do_sample` | Whether Do Sample | Decides whether to use sampling strategies to generate text. If True, employs temperature, top-k, or top-p sampling; if False, uses greedy search or Beam Search. | Set to True for exploratory, creative text generation to increase diversity. Set to False for tasks requiring high precision, such as question answering. |

The main steps of LoRA fine-tuning are as follows:

1. **Model Loading**: Load a **pre-trained** model and tokenizer.

2. **Training**: Train the model with the low-rank matrices while keeping the original model parameters frozen.

3. **Inference**: Generate text based on the input question using the fine-tuned model.

4. **Vote and Evaluation**: Aggregate the predictions from multiple inference paths and select the final answer based on majority voting. Compare the predicted answer with the ground truth to evaluate the model's performance.

**Algorithm 3** Process Questions and Extract Answers

1: **for each** row in df **do**
2:    $question \leftarrow row['question'] +' Options :'$
3:    **for each** option in row['options'].split(',') **do**
4:      $question \leftarrow question+option.replace("]","").replace("[","").replace("'","")$
5:    **end for**
6:    $combined\_question \leftarrow system\_prompt + question$
7:    $max\_attempts \leftarrow 5$
8:    $attempts \leftarrow 0$
9:    $answers \leftarrow []$
10:    $success \leftarrow$ **False**
11:    **while** $\neg success \wedge attempts < max\_attempts$ **do**
12:      **for** $attempt \leftarrow 0$ **to** $2$ **do**
13:        **if** $attempt\%3 == 0$ **then**
14:          $current\_hyperparameters \leftarrow hyperparameters\_1$
15:        **else if** $attempt\%3 == 1$ **then**
16:          $current\_hyperparameters \leftarrow hyperparameters\_2$
17:        **else**
18:          $current\_hyperparameters \leftarrow hyperparameters\_3$
19:        **end if**
20:        $output \leftarrow pipe(combined\_question, **current\_hyperparameters)$
21:        $generated\_text \leftarrow output[0]['generated\_text']$
22:        $correct\_option, success \leftarrow extract\_answer\_final(generated\_text,' A-E')$
23:        **if** $success$ **then**
24:          $answers.append(correct\_option)$
25:        **else**
26:          $attempts \leftarrow attempts + 1$
27:        **end if**
28:      **end for**
29:      **if** $\neg success$ **then**
30:        $df.at[row.index,' y'] \leftarrow "Error"$
31:      **end if**
32:    **end while**
33:    **if** $answers$ **then**
34:      $most\_common\_answer, count \leftarrow Counter(answers).most\_common(1)[0]$
35:      **if** $count > 1$ **then**
36:        $df.at[row.index,' y'] \leftarrow most\_common\_answer$
37:      **else**
38:        $final\_output \leftarrow pipe(combined\_question, **hyperparameters\_2)$
39:        $final\_generated\_text \leftarrow final\_output[0]['generated\_text']$
40:        $final\_correct\_option, final\_success \leftarrow$ $extract\_answer\_final(final\_generated\_text,' A - E')$
41:        **if** $final\_success$ **then**
42:          $df.at[row.index,' y'] \leftarrow final\_correct\_option$
43:        **else**
44:          $df.at[row.index,' y'] \leftarrow "Error"$
45:        **end if**
46:      **end if**
47:    **else**
48:      $df.at[row.index,' y'] \leftarrow "Error"$
49:    **end if**
50: **end for**

## 2.4 Multi-path step-aware inference
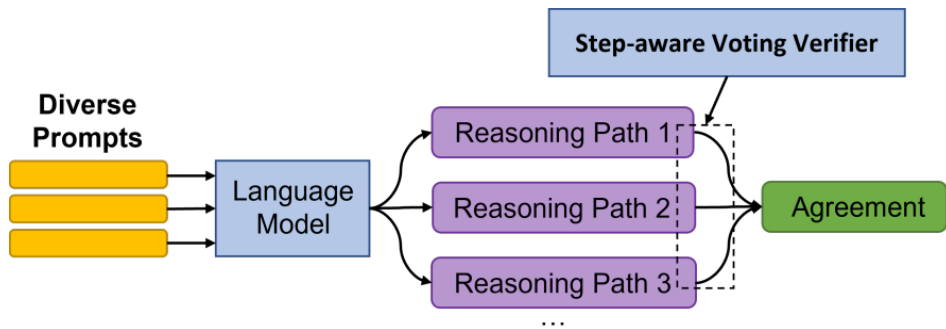
- Specific code can be found in `code/Step-Aware.ipynb`

## 2.4.1 Introduction for DiVeRSe

Prior work (e.g. CoT introduced before) proposed guiding language models with prompts that elicit a series of reasoning steps before giving the final answer. CoT reasoning inserts multi-step paths to improve performance on benchmarks like GSM8K. Few-shot learning tasks require language models to generalize from limited examples. Despite progress, large language models still face difficulties in reasoning tasks such as arithmetic problems. Simple CoT can be regarded as single-vote. In contrast, here we introduce the idea of multi-vote.
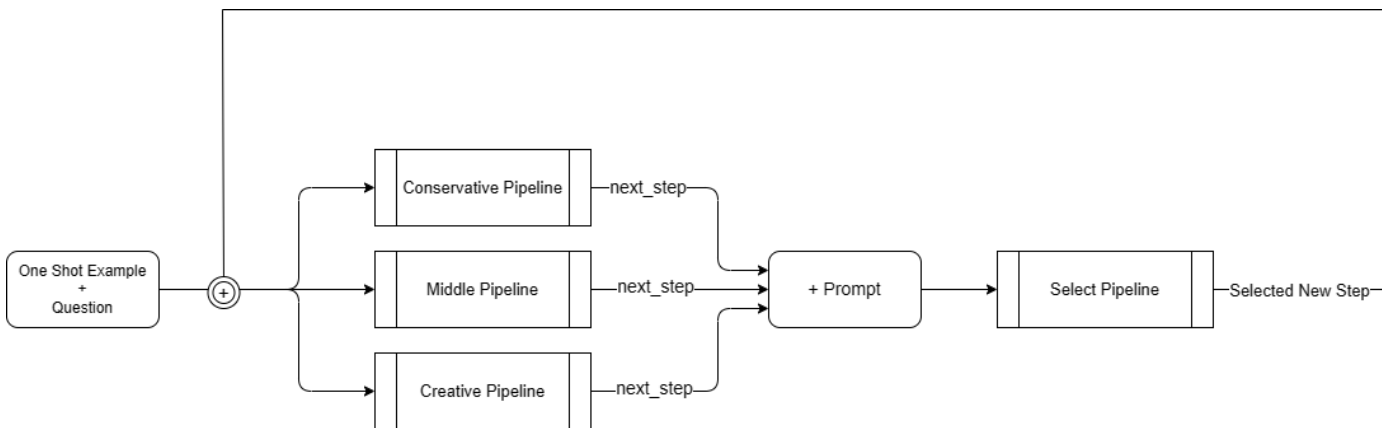
To further improve the performance of CoT, Yifei Li et al proposed a reasoning method called DiVeRSe (Diverse Verifier on Reasoning Step). It has three main components:

- **Diverse Prompts:** Instead of using a fixed set of exemplars, DiVeRSe generates diverse prompts to explore different reasoning paths.
- **Voting Verifier:** Combines voting and a verifier to score reasoning paths and select the most consistent answer.
- **Step-Aware Verifier:** Attributes correctness or wrongness to each reasoning step individually.



With these three processes, DiVeRSe achieves new state-of-the-art results on reasoning benchmarks, including GSM8K (74.4% → 83.2%).

## 2.4.2 Our implementation



Motivated by DiVeRSe, we proposed a multi-path step-aware inference method. The main steps are as follows:

Our method uses multiple path to infer step by step to get the answer of the question.
The method is composed of several reasoning pipelines (RPs) and one select pipeline (SP).

- At each step, all RPs are sent the same prompt(**Combined with a one shot example to cast light on how exactly the COT process should be**), and generate their output respectively.

- Then the SP selects one(Only the **STEP** output but not the whole thing.We do this by splitting the output by the **number tag**) from each output of RPs as the additional prompt in the next round.

- The seletced output will be concatenated to the previous input sentence to become the input of the next round. In a way, this method can be seen as a combination of the previous two methods.

- The input will go through three RPs with different hyperparamaters to infer the next step. In our setting, the three RPs are in conservative, middle and creative mode respectively. After trimming the outputs, the three steps will be inputed into the SP to select one from them as the real next step. (as described above)

- When all the RPs have got their own answer, we will use **multi-vote** to determine the final answer.

## 2.4.3 Example Specification

The following example specifies how it works in the first loop:

Here's the one shot example prepared for the model

```
1   question: What is the least value of x, So that 2x5475 is divisible by 9? options:
    A)7, B)8, C)4, D)3, E)2
2   Let's think step by step.
3   1. The sum of the digits of the number is divisible by 9. Then the number is divisible
    by 9.
4   2. We have 2 + x + 5 + 4 + 7 + 5 = 23 + x.
5   3. The number after 23 and divisible by 9 is 27.
6   4. Least value of x may be '4', so that the total 23 + 4 = 27.
7   5. Answer: C.
```

Combined with the question, we get the initial input.

```
1   one shot example above
2   +
3   question: A man has 12 apples. He eats 4 apples and how many apples does he have now?
    options: A)16, B)12, C)8, D)4, E)2
4   Let's think step by step.
```

Input it into three different RPs to predict next step to solve the question. `hyparamsi` stands for the output of the ith RPs.

```
1   Hyparams0: The man initially has 12 apples. 4 apples are eaten, so we have the formula
    12 - 4 = 8.
2   Hyparams1: The man initially has 12 apples. 4 apples are eaten, so we have the formula
    12 + 4 = 16.
3   Hyparams2: The man eats 4 apples so we should minus 4 from initial number of apples.
```

Input the three possible step into another SP to select one from them. Some necessary prompt and background information are added, as shown below.

```
1   question: A man has 12 apples. He eats 4 apples and how many apples does he have now?
    options: A)16, B)12, C)8, D)4, E)2
2   Given background above, which of following is the next most reasonable inference?
3   #1: The man initially has 12 apples. 4 apples are eaten, so we have the formula 12 - 4
    = 8.
4   #2: The man initially has 12 apples. 4 apples are eaten, so we have the formula 12 + 4
    = 16.
5   #3: The man eats 4 apples so we should minus 4 from initial number of apples.
6   Answer: #
```

The number from SP tells us which step we should choose. Assume we get a "1", then we will choose the first step and add it to the current reason path. The whole reason path now is shown below.

```
1    question: What is the least value of x, So that 2x5475 is divisible by 9? options:
     A)7, B)8, C)4, D)3, E)2
2    Let's think step by step.
3    1. The sum of the digits of the number is divisible by 9. Then the number is
     divisible by 9.
4    2. We have 2 + x + 5 + 4 + 7 + 5 = 23 + x.
5    3. The number after 23 and divisible by 9 is 27.
6    4. Least value of x may be '4', so that the total 23 + 4 = 27.
7    5. Answer: C.
8    question: A man has 12 apples. He eats 4 apples and how many apples does he have
     now? options: A)16, B)12, C)8, D)4, E)2
9    Let's think step by step.
10   1. The man initially has 12 apples. 4 apples are eaten, so we have the formula 12 -
     4 = 8.
```

Now our reasoning has gotten one step. Pass it into RPs again to predict next step and repeat the procedure. If some of the RPs have get an answer, then that RP will stop reasoning.

The step reasoning finishes until the given step number is reached or three PRs have gotten the answer. A mutil vote is used to select one from them as the final answer.

### 2.4.4 Inference Pipeline

---

**Algorithm 4** Inference Procedure

---

1: **procedure** INFERENCE($input, hyper\_list$)
2:     $i \leftarrow 1$
3:     $current\_inf \leftarrow input$
4:     $infs \leftarrow []$
5:     $answers \leftarrow []$
6:     $delete \leftarrow []$
7:     **while** $i < 15 \wedge len(answers) < 3$ **do**
8:         **for** $j \leftarrow 0$ **to** $len(hyper\_list) - 1$ **do**
9:             **if** $j \in delete$ **then**
10:                 **continue**
11:             **end if**
12:             $success \leftarrow$ **True**
13:             $output \leftarrow$ **None**
14:             $inf \leftarrow$ **None**
15:             **if** $success$ **then**
16:                 $output \leftarrow pipe(current\_inf, hyper\_list[j])$
17:                 $inf \leftarrow trim(output[0]['generated\_text'], i)$
18:             **else**
19:                 $delete.append(j)$
20:             **end if**
21:             **if** "Answer" **in** $inf$ **then**
22:                 $ans \leftarrow extract\_answer(inf, True)$
23:                 $answers.append(ans)$
24:                 $delete.append(j)$
25:             **else**
26:                 $infs.append(inf)$
27:             **end if**
28:         **end for**
29:         **if** $len(delete) == 3$ **then**
30:             **break**
31:         **end if**
32:         $background \leftarrow current\_inf.split(shot\_answer)[1]$
33:         $choose \leftarrow select(background, infs)$
34:         **if** $choose == -1 \vee choose \geq len(infs)$ **then**
35:             $choose \leftarrow 0$
36:         **end if**
37:         $current\_inf \leftarrow current\_inf + str(i) + "." + infs[choose]$
38:         $i \leftarrow i + 1$
39:         $infs.clear()$
40:     **end while**
41:     **return** $answers$
42: **end procedure**

---

# 3. Experiment Results

Method of multi-vote is tested on the benchmark AQuA which consists of about 100,000 algebraic word problems. We may just use a subset of it.

## 3.1 Multi-vote

Here are the intermediate inference steps of a quesiton using multi-vote. We set up 3 paths, each of which has a first prompt and a second prompt. Then we select the answer with the highest number of occurrences among 3 paths as the final answer.

```
1    0 Answer: E
2    first_output with prompt0: : The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
3    ...
4    1.  The distance covered by the man in
5
6    second_output with prompt0:  The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
7    ...
8    : A man can row
9
10   first_output with prompt1: : The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
11   ...
12   The speed of the man in still water is 25 kmph. The speed of the man in still water
     is 25 km
13
14   second_output with prompt1:  The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
15   ...
16   The speed of the man in
17
18   first_output with prompt2: : The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
19   ...
20   The mathematical formula to solve this is: 80/25 = 3/11
21
22   second_output with prompt2:  The speed at which a man can row a boat in still water
     is 25 kmph. If he rows downstream, where the speed of current is 11 kmph, what time
     will he take to cover 80 metres?
23   ...
24   The mathematical formula to solve this is: 80/25 = 3/11
25
26   Therefore, among A through E, the answer is C.
27   : A man can row a boat
28   Predicted Answer: A
```

- By the mean of multi-vote, our model can achieve an **accuracy of 20%**.

- While Actually we only tried on 5 questions, the result is not stable enough. But due to the limitation of time and resources, we can't do more experiments. We believe that with more data and more experiments, the model can achieve a higher accuracy.

## Reflections

- **Model Robustness**: The multi-vote method can improve the robustness of the model by exploring different reasoning paths. By aggregating predictions from multiple paths, the model can leverage diverse strategies to enhance its performance.

Having limited time and resources, we can't do more experiments. But there's still a lot of work to do to improve the model's performance.

- We can train the model on a larger dataset beforehead to improve the model's performance.

- We can also try **different** and **more** prompts and hyperparameters to explore more diverse reasoning paths.

## 3.2 LoRA fine-tuning

Here are the intermediate inference steps by which the model after LoRA fine-tuning solves a problem. Like multi-vote, in order to achieve higher stability and reliability, we try to get answers from the model for several times, but here we use the same prompt but different hyperparameters for each time. Then we select the answer with the highest number of occurrences as the final answer.

```
1  Text: You are an expert in Algebra. Lets think through these questions step by step
   to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
   or E. Ensure that your answer follows this format: Answer: . Two friends plan to
   walk along a 43-km trail, starting at opposite ends of the trail at the same time.
   If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
   have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
2  Answer: C)22
3    Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 0.5,
   'top_k': 50, 'top_p': 0.95, 'num_return_sequences': 1, 'pad_token_id': 2}
4
5  Question ID 0, Attempt 0: Correct Option - C
6
7  Text: You are an expert in Algebra. Lets think through these questions step by step
   to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
   or E. Ensure that your answer follows this format: Answer: . Two friends plan to
   walk along a 43-km trail, starting at opposite ends of the trail at the same time.
   If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
   have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
8  Answer: C)22
9    Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 0.7,
   'top_k': 80, 'top_p': 0.9, 'num_return_sequences': 1, 'pad_token_id': 2}
10
11 Question ID 0, Attempt 1: Correct Option - C
12
```

```
13   Text: You are an expert in Algebra. Lets think through these questions step by step
     to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
     or E. Ensure that your answer follows this format: Answer: . Two friends plan to
     walk along a 43-km trail, starting at opposite ends of the trail at the same time.
     If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
     have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
     Answer:. An electron is accelerated by a voltage of 500 V in a region where the
     electric field is 150 N/C. If the electron's mass is 9.11 × 10⁻² kg, what is its
     acceleration?['A)10.0 N/C', 'B)10.1 N/C', 'C)10.2 N/C', 'D)10.
14     Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 1.0,
     'top_k': 100, 'top_p': 0.85, 'num_return_sequences': 1, 'pad_token_id': 2}
15   Attempt 2 failed, retrying...
16
17   Text: You are an expert in Algebra. Lets think through these questions step by step
     to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
     or E. Ensure that your answer follows this format: Answer: . Two friends plan to
     walk along a 43-km trail, starting at opposite ends of the trail at the same time.
     If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
     have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
18   Answer:. Two friends plan to walk along a 43-km trail, starting at opposite ends of
     the trail at the same time. If Friend P's rate is 15% faster than Friend Q's, how
     many kilometers will Friend P have walked when they pass each other?['A)21',
     'B)21.5', 'C)22', 'D)22.5', 'E)23']
19   An
20     Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 0.5,
     'top_k': 50, 'top_p': 0.95, 'num_return_sequences': 1, 'pad_token_id': 2}
21   Attempt 3 failed, retrying...
22
23   Text: You are an expert in Algebra. Lets think through these questions step by step
     to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
     or E. Ensure that your answer follows this format: Answer: . Two friends plan to
     walk along a 43-km trail, starting at opposite ends of the trail at the same time.
     If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
     have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
     Answer:. A company produces 20,000 units of a product in January, 15,000 units in
     February, and 25,000 units in March. What is the total production for the three
     months?['A)60,000', 'B)61,500', 'C)62,000', 'D)62,500', 'E)63,0
24     Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 0.7,
     'top_k': 80, 'top_p': 0.9, 'num_return_sequences': 1, 'pad_token_id': 2}
25   Attempt 4 failed, retrying...
26
27   Text: You are an expert in Algebra. Lets think through these questions step by step
     to ensure we get the correct answer. There are only 5 options to choose from A,B,C,D
     or E. Ensure that your answer follows this format: Answer: . Two friends plan to
     walk along a 43-km trail, starting at opposite ends of the trail at the same time.
     If Friend P's rate is 15% faster than Friend Q's, how many kilometers will Friend P
     have walked when they pass each other?['A)21', 'B)21.5', 'C)22', 'D)22.5', 'E)23']
28   Answer: E)23
29     Hyperparameters: {'do_sample': True, 'max_new_tokens': 100, 'temperature': 1.0,
     'top_k': 100, 'top_p': 0.85, 'num_return_sequences': 1, 'pad_token_id': 2}
30
31   Question ID 0, Attempt 5: Correct Option - E
```

```
32
33  Question ID 0: Majority Option - C
```

- By the mean of LoRA fine-tuning, our model can achieve an **accuracy of 26%**.(On 100 questions)

**Possible improvement:**

- In our current model, we only use 3 hyperparameters for inference. We can try more hyperparameters to explore more diverse reasoning paths and improve the model's performance.

- We can enlarge `max_attempts` to get more stable results.

## 3.3 Multi-path

Here is an example of a question, its true answer and the predicted answer given by the model using multi-path.

```
1  Question: There are k-2 members in a certain band, including Jim and Ellen. Two
   members are to be selected to attend the Grammy awards ceremony. If there are 6
   possible combinations in which Jim and Ellen are not selected, what is the value of k?
2  A)8
3  B)9
4  C)10
5  D)11
6  E)12
7  True Answer: A
8  Predicted Answer: A
```

Here are the intermediate inference steps of this quesiton using multi-path. We set up 3 groups of hyperparameters, which represent 3 reasoning pipelines respectively. The model does inference using these 3 groups of hyperparameters parallelly, and then gets 3 outputs. Then the select pipeline chooses the most reasonable outputs among them as the additional prompt in the next round. Finally the answer with the highest number of occurrences is selected as the final answer.

```
1   question 0 : correct answer is A
2   Hyperams0:The band has k-2 members.
3   Hyperams1:Total members in the band = k.
4   Hyperams2:Number of members in the band, including Jim and Ellen, is k.
5   Hyperams0:Jim and Ellen are in the band.
6   Hyperams1:Jim and Ellen are in the band.
7   Hyperams2:There are 6 possible combinations of 2 members in the band that do not
    include Jim and Ellen.
8   Hyperams0:There are 6 possible combinations in which Jim and Ellen are not selected.
9   You seem to be using the pipelines sequentially on GPU. In order to maximize
    efficiency please use a dataset
10  Hyperams1:Two members are to be selected for the Grammy awards.
11  Hyperams2:There are 6 possible combinations in which Jim and Ellen are not selected.
12  Hyperams0:Let's call the number of combinations where Jim and Ellen are not selected
    as x.
13  Hyperams1:The number of combinations in which Jim and Ellen are not selected is 6.
14  Hyperams2:So the total number of possible combinations is (k-2) \* 6 = k \* 6.
15  Hyperams0:The total number of combinations is k - 1.
```

```
16  Hyparams1:Since Jim and Ellen are in the band, the total number of combinations is
    k.
17  Hyparams2:Then, the number of combinations in which Jim and Ellen are selected is x
    + 1.
18  Hyparams0:The number of combinations where Jim and Ellen are selected is k - x.
19  Hyparams1:x + 1 = k - 1, so x = k - 2.
20  Hyparams2:Jim and Ellen are selected from the remaining combinations.
21  Hyparams0:The number of combinations is x.
22  Hyparams1:The value of k can be calculated as k = 6 + x = 6 + (k - 2) = 8.
23  Hyparams2:The total number of combinations where Jim and Ellen are not selected is x
    + 1 = (k - 1) - 2 = k - 3.
24  Hyparams0:The total number of combinations is k - 1 + x.
25  Hyparams1:The total number of combinations is k - 1.
26  Hyparams2:The total number of combinations is (k - 1) x.
27  Hyparams0:The number of combinations where Jim and Ellen are not selected is x.
28  Hyparams1:The number of combinations where Jim and Ellen are not selected is x.
29  Hyparams2:Since Jim and Ellen are in the band, the number of combinations is k - 1 +
    k - 2 = k.
30  Hyparams0:The total number of combinations is k - 1 + x = k - 1 + x = (k - 1) + x =
    k.
31  Hyparams1:The number of combinations where Jim and Ellen are selected is x + 1.
32  Hyparams2:The total number of combinations where Jim and Ellen are not selected is k
    - 1 + x.
33  Hyparams0:Answer: A.
34  Hyparams1:The number of combinations is x.
35  Hyparams2:Since Jim and Ellen are in the band, the total number of combinations is
    k.
36  Hyparams1:The total number of combinations is k.
37  Hyparams2:x + 1 = k - 1, so x = k - 2.
38  Hyparams1:The number of combinations where Jim and Ellen are not selected is x.
39  Hyparams2:The number of combinations where Jim and Ellen are not selected is x.
40  Hyparams1:The total number of combinations is k.
41  Hyparams2:The total number of combinations is k.
42  answer: A
43  Predicted Answer: A
```

- By the mean of multi-path, our model can achieve an **accuracy of 21%**(on **40** questions).

**Possible improvement:**

- Look at the output below,we can find that the inference was stopped before it reached all the answers. We can enlarge the number of steps to get more stable results.

```
1   correct answer is D
2   Hyparams0:k/l < 1, so k and l are both greater than 1.
3   Hyparams1:If k/l < 1, then k and l are both less than 1.
4   Hyparams2:If k/l < 1, then k/l is less than 1, and therefore k is greater than l.
5   Hyparams0:k/l < 1, so k^2/l^2 < 1.
6   Hyparams1:Both k and l are positive integers, so k^2 and l^2 are also positive
    integers.
7   Hyparams2:The expression k/l^2 is always greater than 1, because the square of any
    positive number is always greater than the number itself.
8   Hyparams0:k^2/l^2 < 1, so k^2/l < 1.
9   Hyparams1:k^2/l^2 < 1, so k^2/l < 1.
10  Hyparams2:C)k^2/l < 1.
11  Hyparams0:k^2/l < 1, so k^2 < l.
12  Hyparams1:k^2/l < 1, so k^2 < l.
13  Hyparams2:We know that both k and l are positive integers.
14  Hyparams0:k^2 < l, so k < l.
15  Hyparams1:Since k is positive and greater than 1, k^2 must be greater than 1.
16  Hyparams2:k^2 < l, so l > k^2.
17  Hyparams0:Answer: B.
18  Hyparams1:Let's take C. k^2/l > 1.
19  Hyparams2:k < l, so l is the greatest number.
20  Hyparams1:But k^2/l^2 = k/l < 1, which is not possible.
21  Hyparams2:C. Answer: C.
22  Hyparams1:So, C is not the correct answer.
23  Hyparams1:D. l/k > 1.
24  Hyparams1:k/l > 1, so k > 1.
25  Hyparams1:But k^2/l > 1, so k > l.
26  answer: C
27  Predicted Answer: C
28  correct answer is C
```

- We can also try more hyperparameters to explore more diverse reasoning paths and improve the model's performance.

- We can enlarge the number of paths to get more stable results.

- We can also use the fine tuning method beforehead to improve the model's performance.

# References

[1] Li, Yifei, et al. "Making large language models better reasoners with step-aware verifier." arXiv preprint arXiv:2206.02336 (2022).

[2] Alpaca-CoT https://github.com/PhoebusSi/Alpaca-CoT