# Long Math Reasoning Problem Generation

Changwei Li, Guangping Huang, **Zihao Zhou**, and Qiufeng Wang

XJTLU & UoL

# Why long math reasoning problem generation?

Long context reasoning is important in multiple real-world scenarios



**Paper Interpretation**

Understanding complex research papers requires processing entire entire documents, not just snippets.

**Meeting Minutes**

Summarize hours-long riours=ions discussions, preserving context, decisions, and action items.

**Medical Q,A**

Analyzing patient history, multiple reports, and medical literature for accurate diangoses

**Legal Q,A**

Revieiwg contracts, case law, and regulations to provide informed legal advice

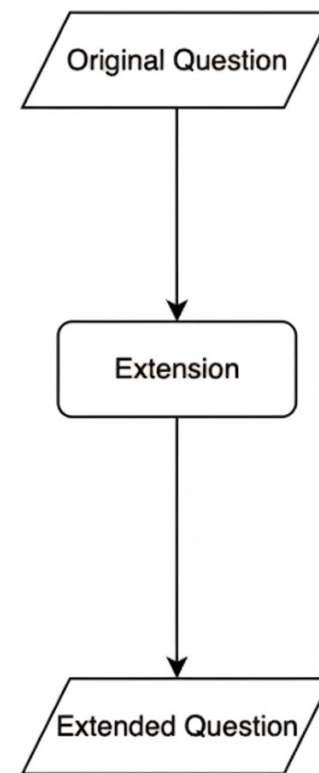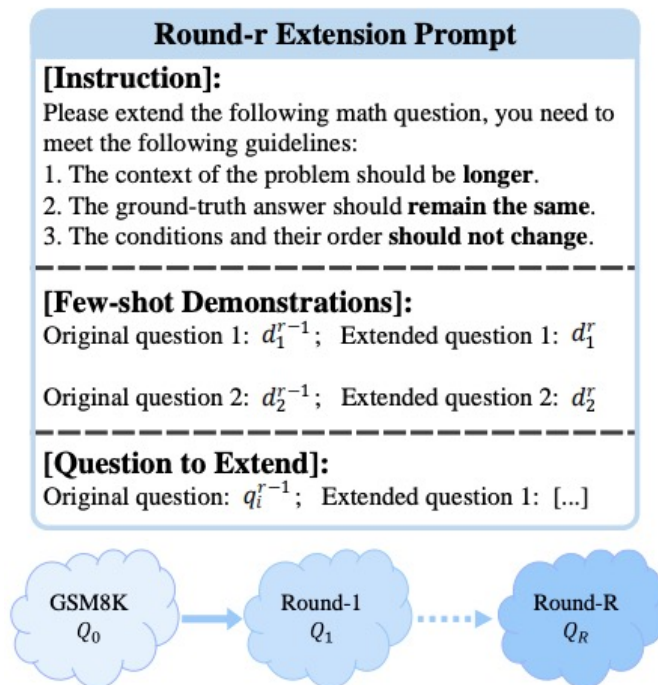# Why long math reasoning problem generation?

Many efforts work on collect or synthesize long reasoning data

  - Evaluation: LongBench, Bamboo, LooGLE, L-evel (<span style="color:orange">acl24 outstanding paper</span>),MiniLongbench(<span style="color:orange">acl25 outstanding paper</span>)

  - Training: AI companies synthesises long- context reasoning data to enhance the foundation model's ability.

Math reasoning problem is high-quality testbed and training data in reasoning tasks, and it is easy to verify too.

In this paper, we propose generating long math reasoning data from existing short math word problems !

- Generating long math reasoning data from scratch is difficult. (long-context, correctness of logic and answer )

- Short math word problems is a good resource as the seed data (reasoning scenario, detailed solution, Label)

**Round-r Extension Prompt**

[Instruction]:
Please extend the following math question, you need to meet the following guidelines:
1. The context of the problem should be **longer**.
2. The ground-truth answer should **remain the same**.
3. The conditions and their order **should not change**.

- - - - - - - - - - - - - - - - - - - - - - - - - - -

[Few-shot Demonstrations]:
Original question 1: $d_1^{r-1}$; Extended question 1: $d_1^r$

Original question 2: $d_2^{r-1}$; Extended question 2: $d_2^r$

- - - - - - - - - - - - - - - - - - - - - - - - - - -

[Question to Extend]:
Original question: $q_i^{r-1}$; Extended question 1: [...]

GSM8K
$Q_0$

Round-1
$Q_1$

Round-R
$Q_R$

Original Question

↓

Extension

↓

Extended Question

Applying LLMs to extend the short math word problems

Can LLMs Solve Longer Math Word Problems Better? (ICLR 25)

GSM8K (Q0) → Extension → Q1
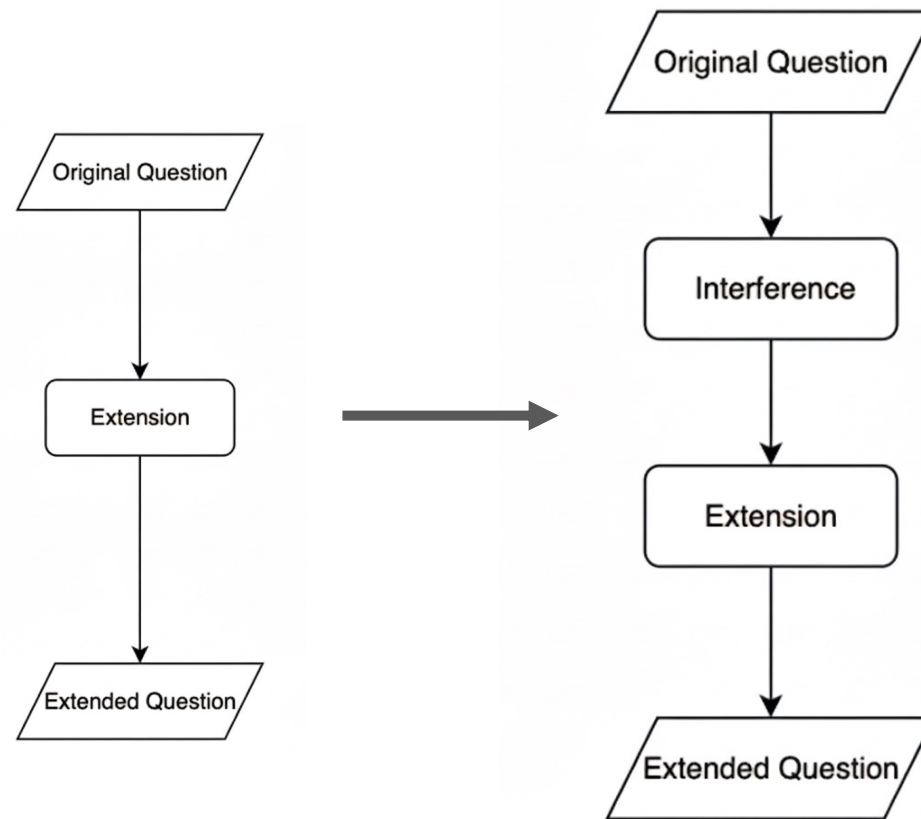
Q1 → Extension → Q2

Q2 → Extension → Q3

. . . . . . .

then repeat this loop to generate longer MWPs

It's too simple and has many limitation…

- The length will no longer increases after a certain number of epochs.

- These data lack noise, whereas real long-context data usually contain noisy information that interferes with reasoning.
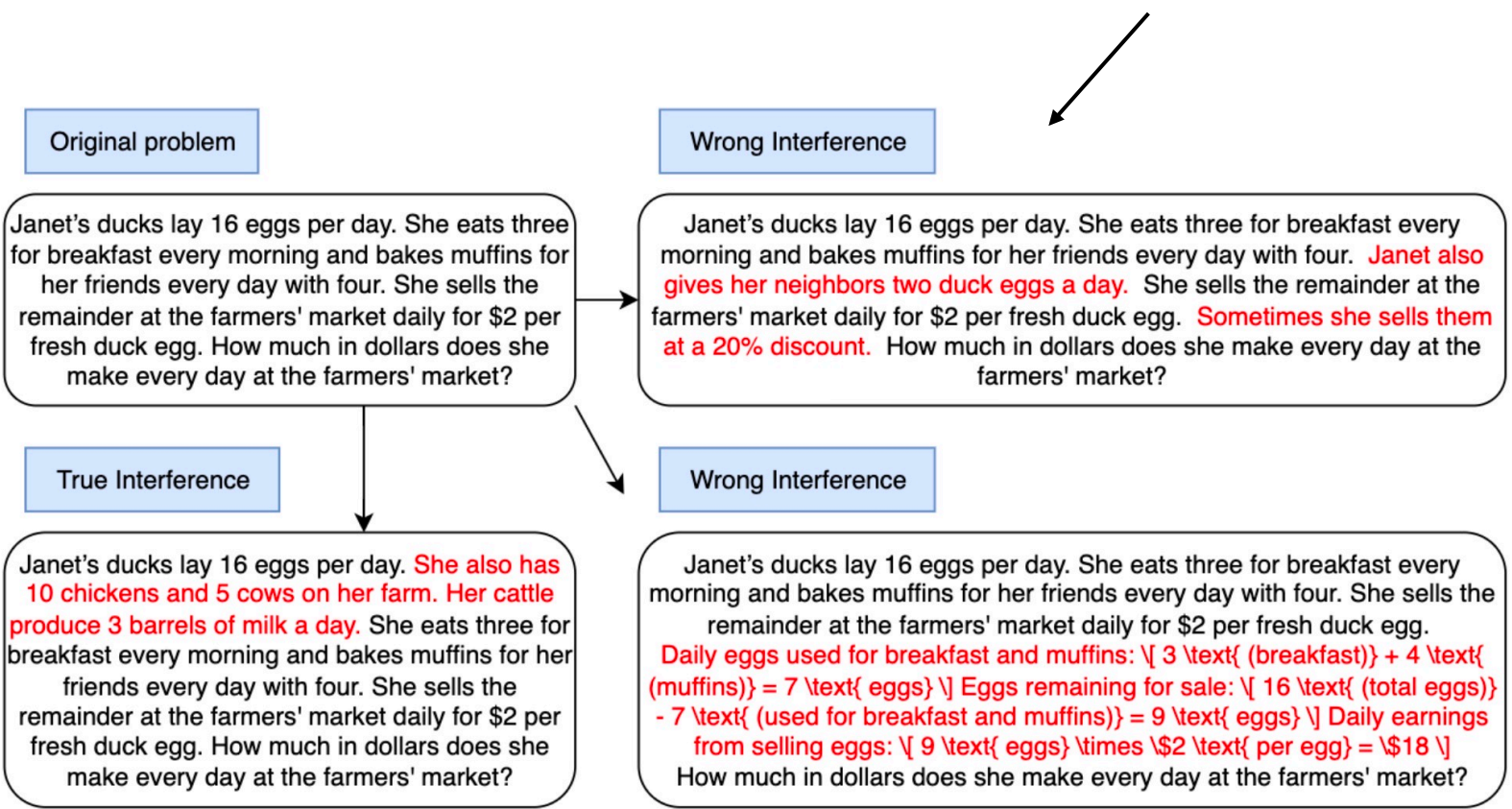
**Original problem**

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?

**True Interference**

Janet's ducks lay 16 eggs per day. She also has 10 chickens and 5 cows on her farm. Her cattle produce 3 barrels of milk a 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?

Original Question → Extension → Extended Question

→

Original Question → Interference → Extension → Extended Question

Add some noise information by LLMs before extension

But it has some problems…

Change the ground truth!

Original problem

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?

Wrong Interference

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. Janet also gives her neighbors two duck eggs a day. She sells the remainder at the farmers' market daily for $2 per fresh duck egg. Sometimes she sells them at a 20% discount. How much in dollars does she make every day at the farmers' market?

True Interference

Janet's ducks lay 16 eggs per day. She also has 10 chickens and 5 cows on her farm. Her cattle produce 3 barrels of milk a 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?

Wrong Interference

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. Daily eggs used for breakfast and muffins: \[ 3 \text{ (breakfast)} + 4 \text{ (muffins)} = 7 \text{ eggs} \] Eggs remaining for sale: \[ 16 \text{ (total eggs)} - 7 \text{ (used for breakfast and muffins)} = 9 \text{ eggs} \] Daily earnings from selling eggs: \[ 9 \text{ eggs} \times \$2 \text{ per egg} = \$18 \] How much in dollars does she make every day at the farmers' market?
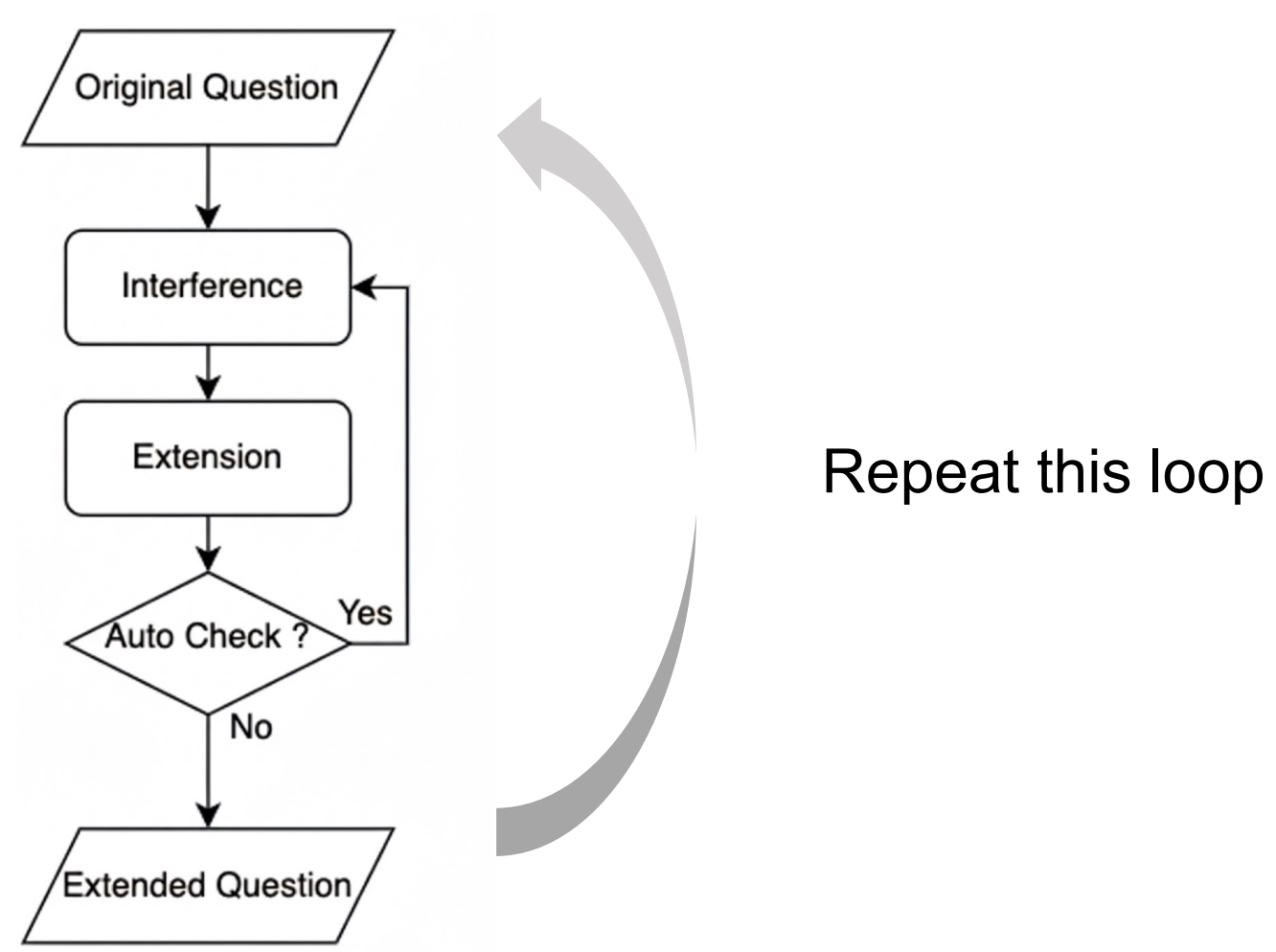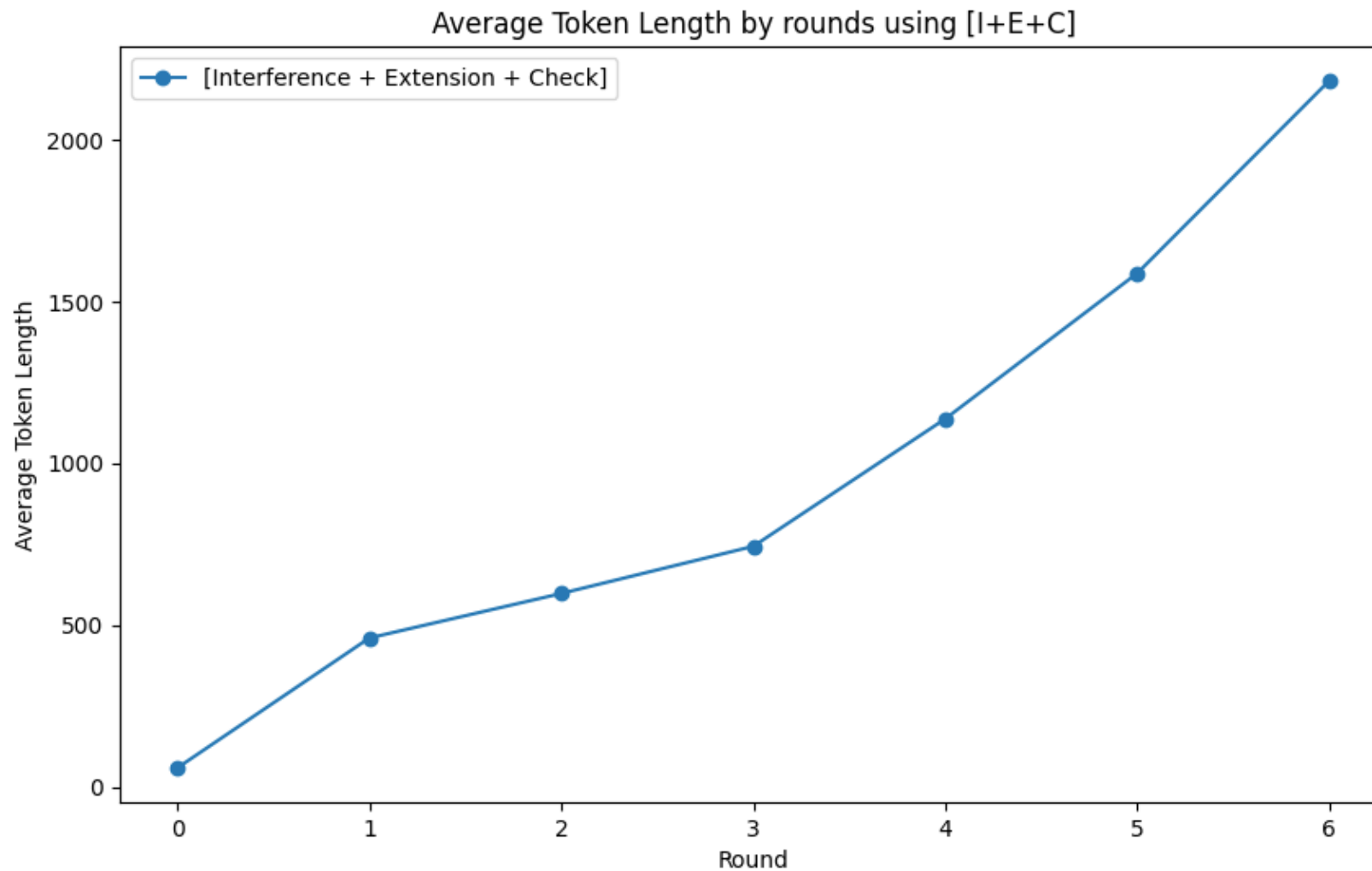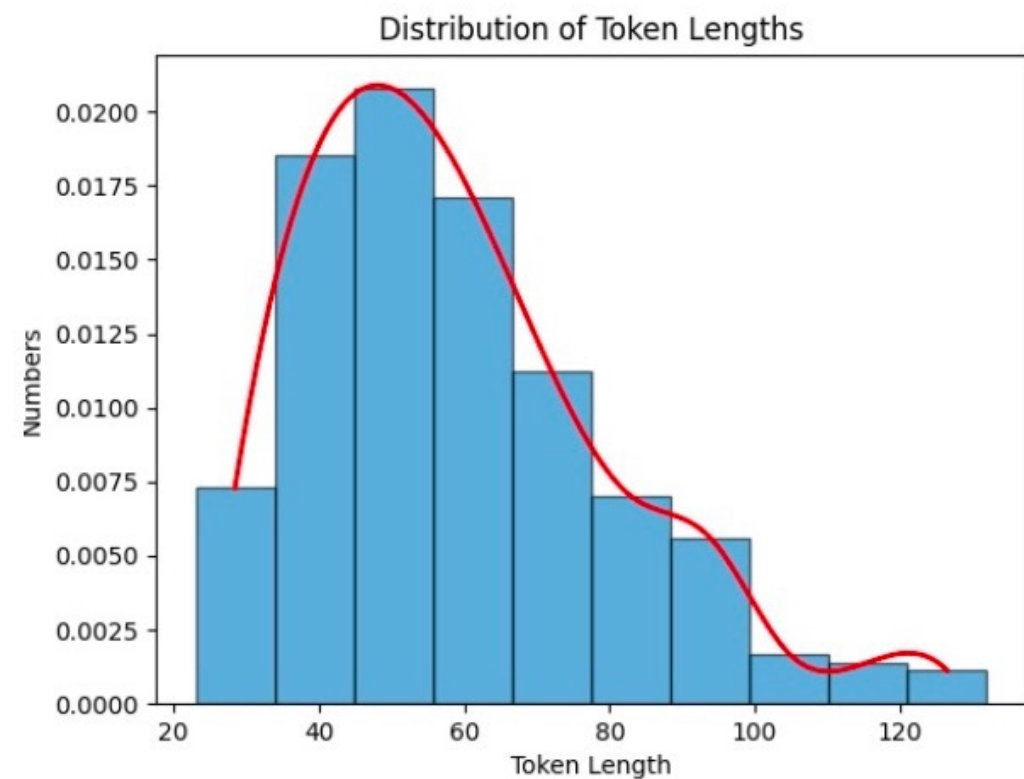
Generate the solution!

LLM-as-a-Judge:
*(1) whether solution steps have been added to the problem, and (2) whether the added distractors have affected the problem's logic and final answer. If issues are found, the problem is returned to the previous round for interference and extension.*
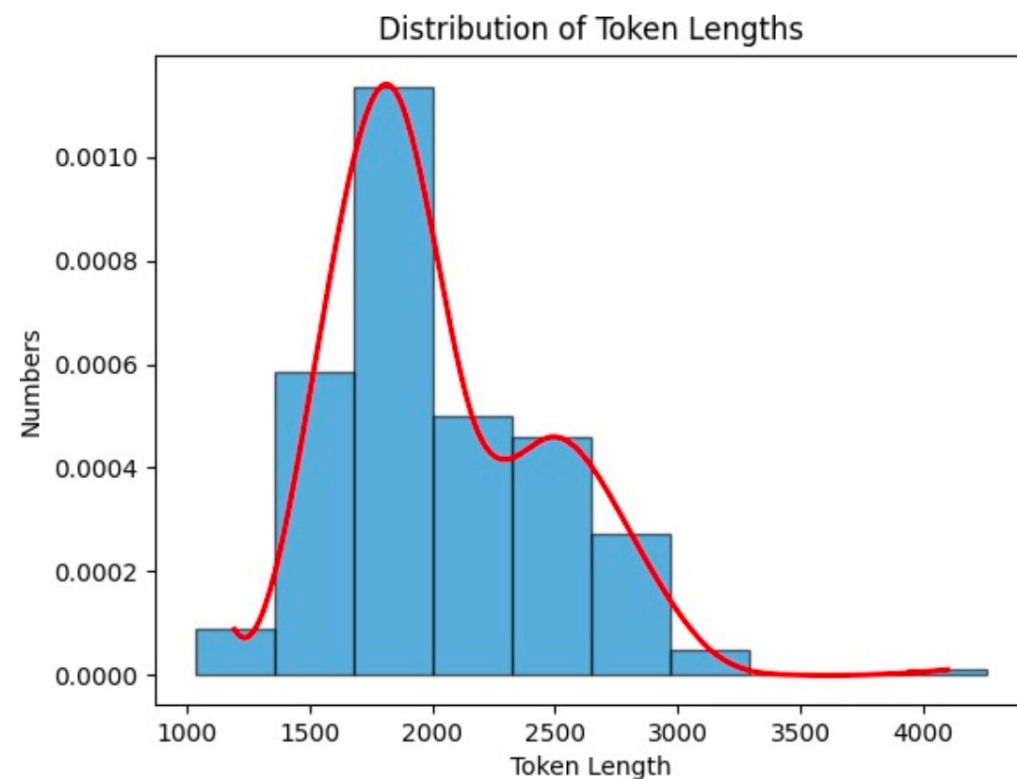
# Final Pipeline:



Repeat this loop
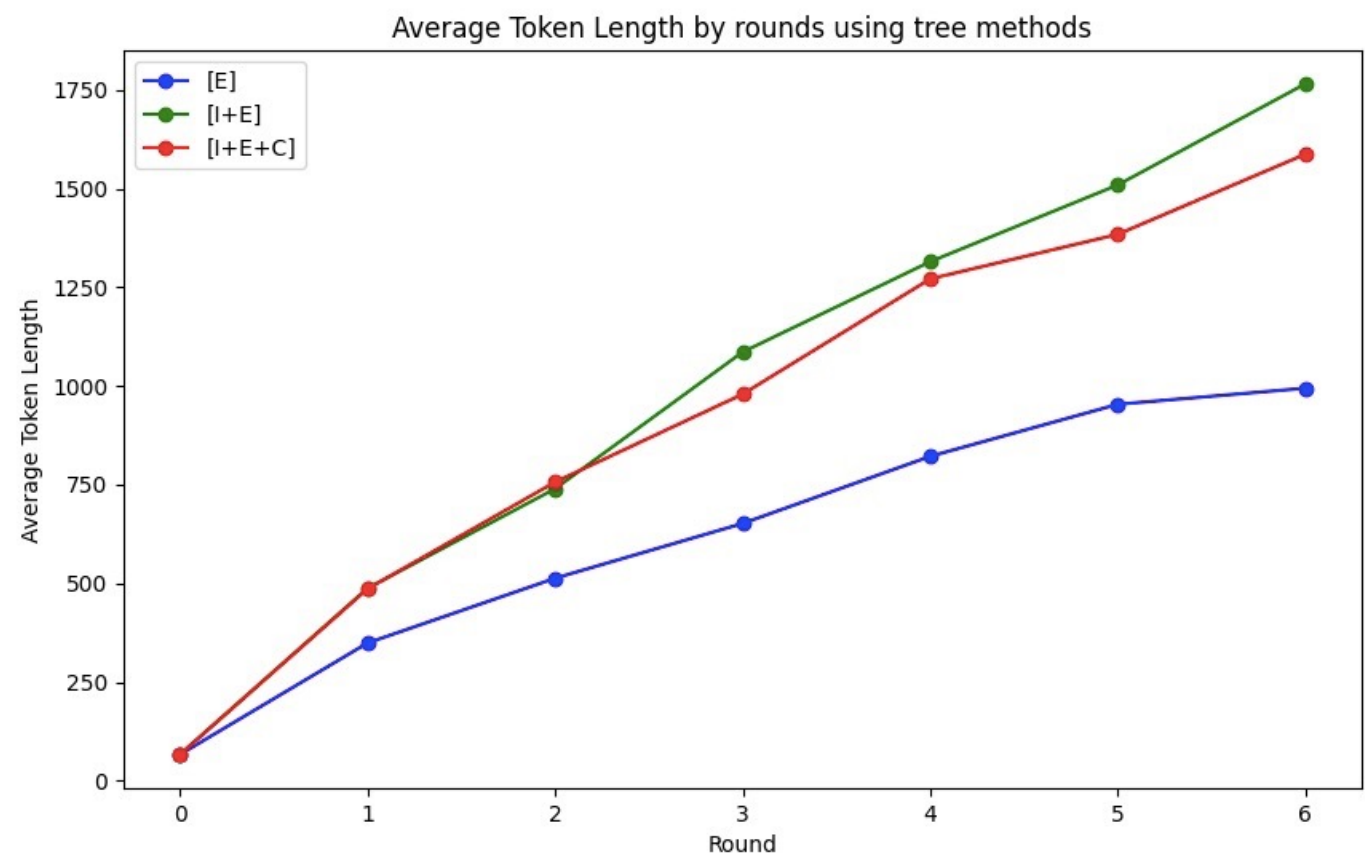
# Generation on GSM8k



Average Token Length by rounds using [I+E+C]

(a) Length in the GSM8K dataset

(b) MWP length after six extensions

# Generation Acc



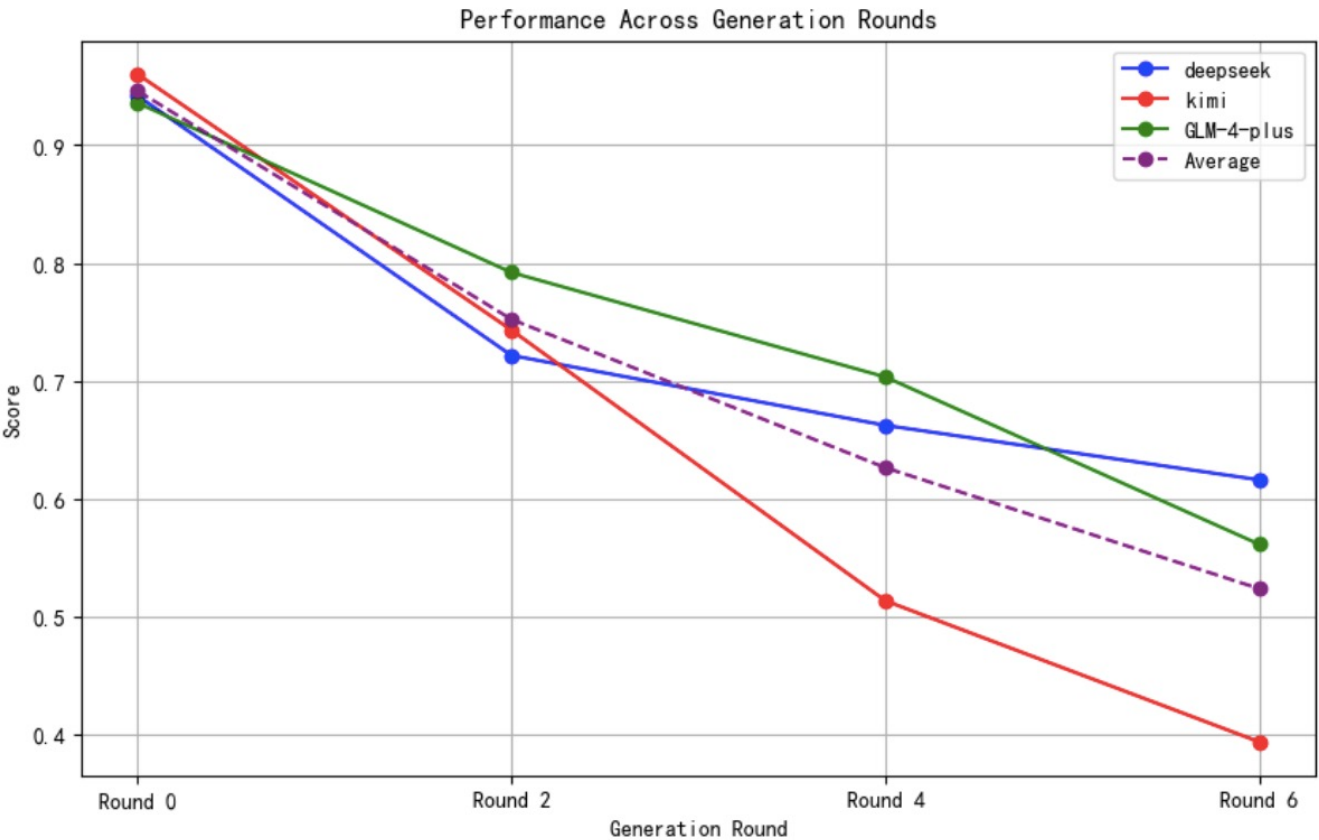| $I + E$ | $I + E + C$ |
| --- | --- |
| 60% | 85% |

# Evalution Acc

## Table 2: The accuracy of the model's answer in different rounds

| Generation round | Round 0 (Original) | Round 2 | Round 4 | Round 6 |
|---|---|---|---|---|
| DeepSeek | 0.9419 | 0.7217 | 0.6624 | 0.6162 |
| Kimi | 0.9602 | 0.7432 | 0.5137 | 0.3940 |
| GLM | 0.9358 | 0.7920 | 0.7034 | 0.5615 |
| Average | 0.9459 | 0.7523 | 0.6265 | 0.5239 |



Performance Across Generation Rounds

# Discussion

Generating long math reasoning data from existing short math word problems is promising

Generating long math reasoning data from environment? (collect action, final signal)

# Thank You!