

# Better Zero-Shot Reasoning with Self-Adaptive Prompting

Muning Wen

10.18

---

# Better Zero-Shot Reasoning with Self-Adaptive Prompting

**Xingchen Wan<sup>\*1,3</sup>, Ruoxi Sun<sup>1</sup>, Hanjun Dai<sup>2</sup>, Sercan Ö. Arık<sup>1</sup>, Tomas Pfister<sup>1</sup>**

<sup>1</sup>Google Cloud AI Research

<sup>2</sup>Google DeepMind

<sup>3</sup>Department of Engineering Science, University of Oxford  
{xingchenw, ruoxis, hadai, soarik, tpfister}@google.com

<https://arxiv.org/abs/2305.14106>

# Motivation

---

## Limitation of zero-shot reasoning:

Its performance is limited due to the lack of guidance to the LLMs.

## Limitation of in-context reasoning:

1. Performance is sensitive to the choice of examples.
2. Designing examples requires significant human effort.
3. The diversity of downstream tasks of LLMs/novel test-time tasks unseen previously.

## The focus of this work:

Improving LLM reasoning ability in the **general zero-shot** setup **with access to input queries but not labels**.

## Key idea:

Collect a pool of rationales and answers to a set of questions with Zero-shot CoT, then **select** the most suitable questions as in-context examples.

## Chain-of-Thought Prompting

### Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓

# Problem Analysis - How impactful are the contextual examples?

## Problem Settings:

1. Zero-shot reasoning (self-adaptive prompting)
2. A set of input queries are available (diverse).
3. Do not need labels.

## The influence of in-context examples:

1. Zero-shot CoT with no demo: correct logic but wrong answer;
2. Correct demo and correct answer;
3. Correct but repetitive demo leads to repetitive outputs;
4. Erroneous demo leads to a wrong answer;
5. Combining erroneous and correct demo leads to a correct answer.

## Thus:

In-context demos need **carefully-designed selection procedure** (key objective of this paper).

**[Question]** Henry had 11 dollars. For his birthday he got 18 more dollars but spent 10 on a new game. How much money does he have now?

**[Demo1]** Q: John bought 21 games from a friend and bought 8 more at a garage sale. If 23 of the games didn't work, how many good games did he end up with? A: Let's think step by step. He bought  $21 + 8 = 29$  games in total. He has  $29 - 23 = 6$  good games. ✓

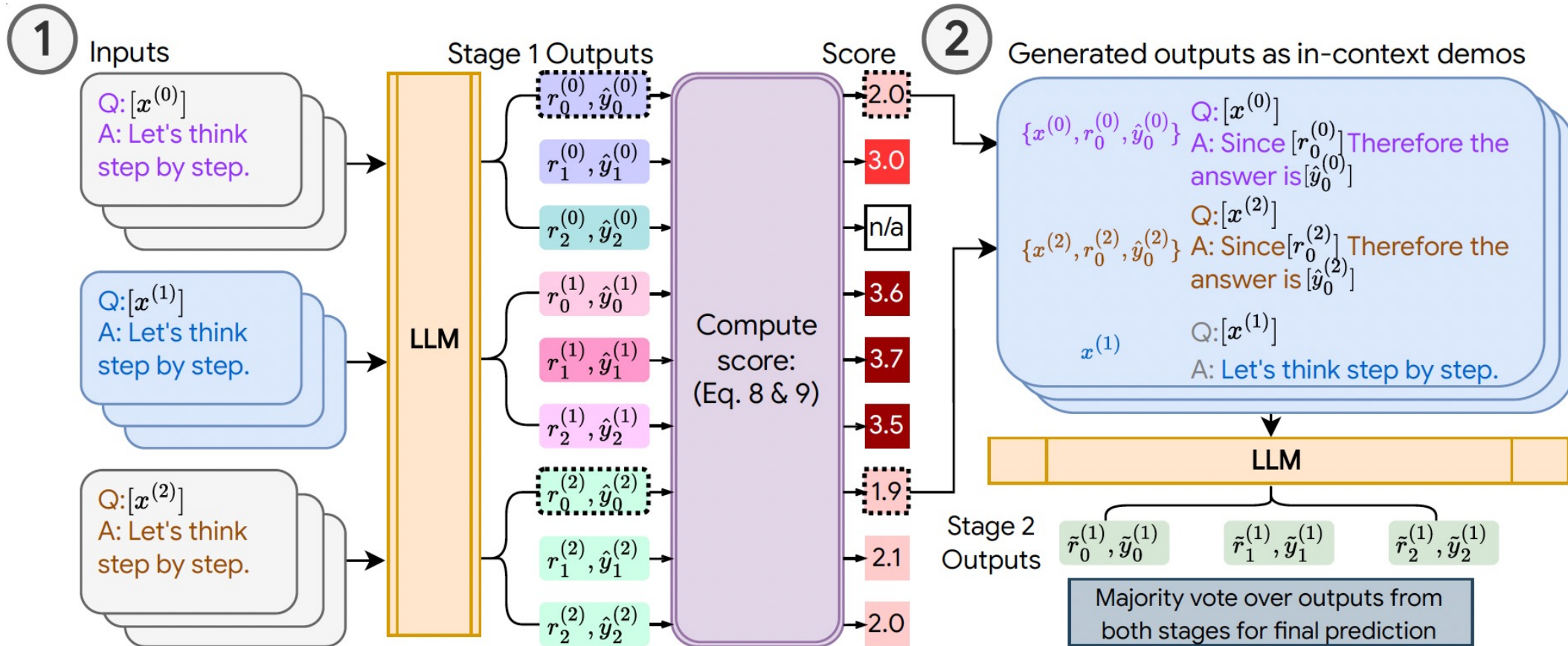
**[Demo2]** Q: John bought 21 games from a friend and bought 8 more at a garage sale. If 23 of the games didn't work, how many good games did he end up with? A: Let's think step by step. He bought  $21 + 8 = 29$  games in total. He bought  $21 + 8 = 29$  games in total. He has  $29 - 23 = 6$  good games. He has  $29 - 23 = 6$  good games. ✓

**[Demo3]** Q: Faye picked 23 and her mother picked 5 carrots from the garden. If only 12 of the carrots were good, how many carrots were bad? A: Let's think step by step.  $23 + 5 = 28$   $28 - 12 = 16$   $16 / 2 = 8$   $8 / 2 = 4$   $4 / 2 = 2$   $2 / 2 = 1$   $1 / 2 = 0$ . ✗

Input	Output
Q: <b>[Question]</b> A: Let's think step by step.	How much money does he have now? Henry has $11 + 18 - 10 = 27$ . ✗
<b>[Demo1]</b> Q: <b>[Question]</b> A: Let's think step by step.	He had $11 + 18 = 29$ dollars. He spent 10 dollars on a game, so he has $29 - 10 = 19$ dollars. ✓
<b>[Demo2]</b> Q: <b>[Question]</b> A: Let's think step by step.	He had 11 dollars. He had 11 dollars. He had 11 dollars. He had 11 dollars.... ✗
<b>[Demo3]</b> Q: <b>[Question]</b> A: Let's think step by step.	$11 + 18 = 29$ $29 - 10 = 19$ $19 / 2 = 9$ $9 / 2 = 4$ $4 / 2 = 2$ $2 / 2 = 1$ $1 / 2 = 0$ ✗
<b>[Demo3] [Demo1]</b> Q: <b>[Question]</b> A: Let's think step by step.	$11 + 18 = 29$ . $29 - 10 = 19$ . ✓



# Framework - Consistency-based Self-adaptive Prompting (COSAP)



## Stage 1:

1. Collects the pool of rationales and answers to question set via Zero-shot CoT.
2. Compute the *score* of each question, a metric inspired by *self-consistency*.
3. Identify suitable question-reasoning pairs with majority vote on the score.

## Stage 2:

1. Augment the target question with a number of selected in-context demonstrations.
2. The augmented question is used to query the LLM a second time.
3. A majority vote over outputs from both stages forms the final prediction.

# Demonstration Selection - How to calculate the score?

## Building the Candidate Pool:

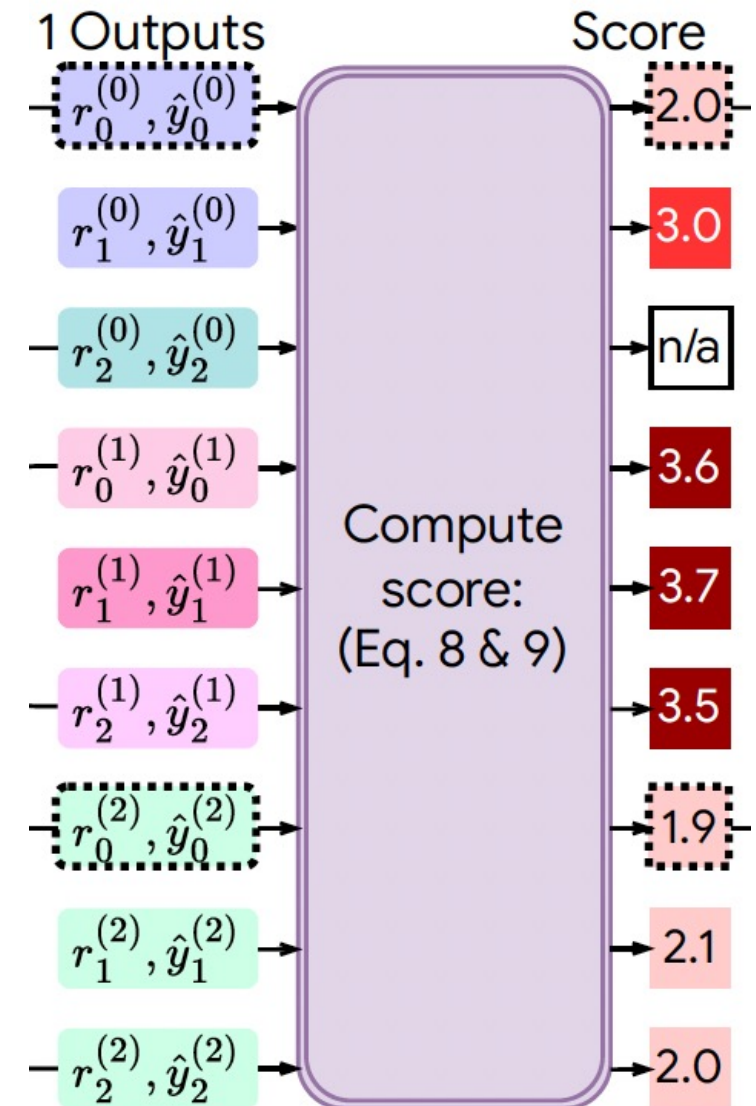
1. Run Zero-shot CoT over all questions.
2. Query the LLM  $m$  times with non-zero temperature for each question.
3. Extract reasoning paths  $\{r_j^{(i)}\}_{j=1}^m$ , and potential answers  $\{\hat{y}_j^{(i)}\}_{j=1}^m$ .

## Demonstration selection is the key objective of this paper:

1. In-context learning is sensitive to the choices of the demonstrations.
2. Select small  $K$  (typically  $\leq 10$ ) demos from a large set of candidates.
3. The candidate pool is imperfect (due to the absence of labels).

## Criteria:

1. Consistency
2. Diversity
3. Repetition



# Self-Consistency (Qualitative)

## Reasons:

1. “Majority predictions are more likely to be correct”.
2. To prune the candidate pool.
3. To select the demonstrations in absence of ground-truths.

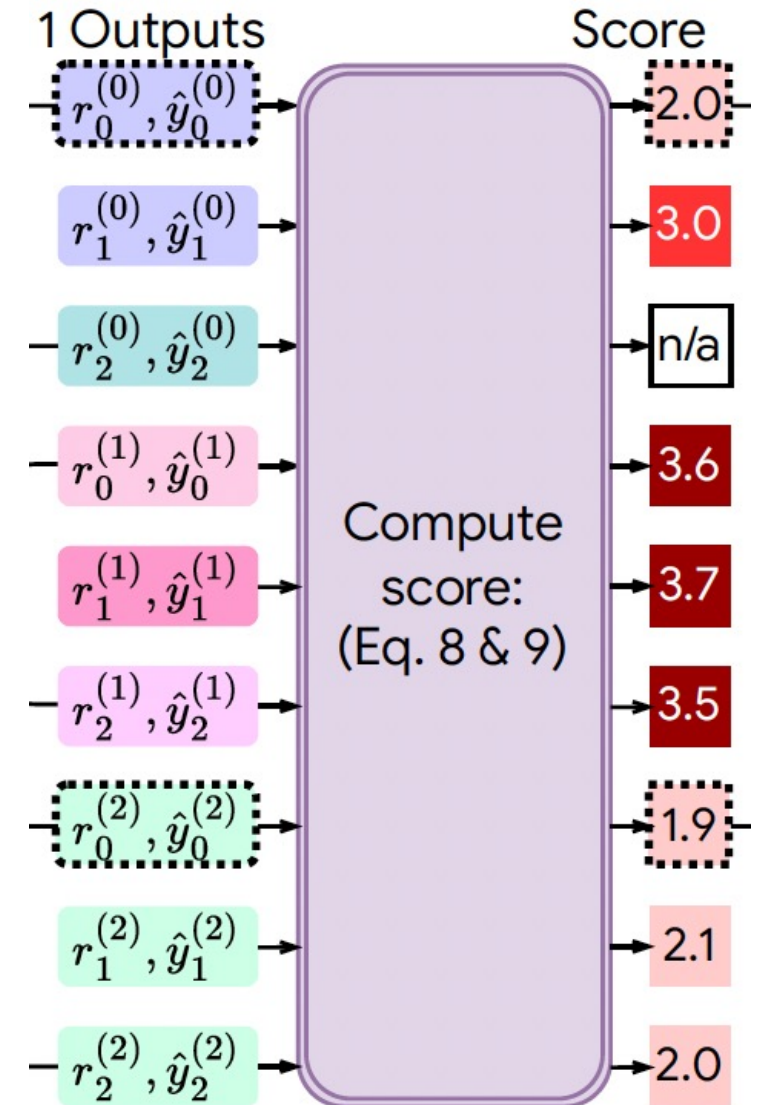
## Measurement:

1. Compute majority vote prediction(s) from all predictions with:

$$\{(r_j^{(i)}, \hat{y}_j^{(i)})\}_{j=1}^m \sim \mathbb{P}(r^{(i)}, \hat{y}^{(i)} | x^{(i)}, c, \theta), \quad (1)$$

$$\hat{y}^{(i)} = \arg \max_{\hat{y}_j^{(i)}} \sum_{k=1}^m \mathbb{I}(\hat{y}_j^{(i)} = \hat{y}_k^{(i)}), \quad (2)$$

2. Retain only the rationales that lead to the majority vote prediction.
3. Use further heuristics to remove obviously bad candidates (e.g. responses containing no numbers for arithmetic tasks, or overly short and/or fragmented responses).



# Self-Consistency (Quantitative)

## Reasons (cont):

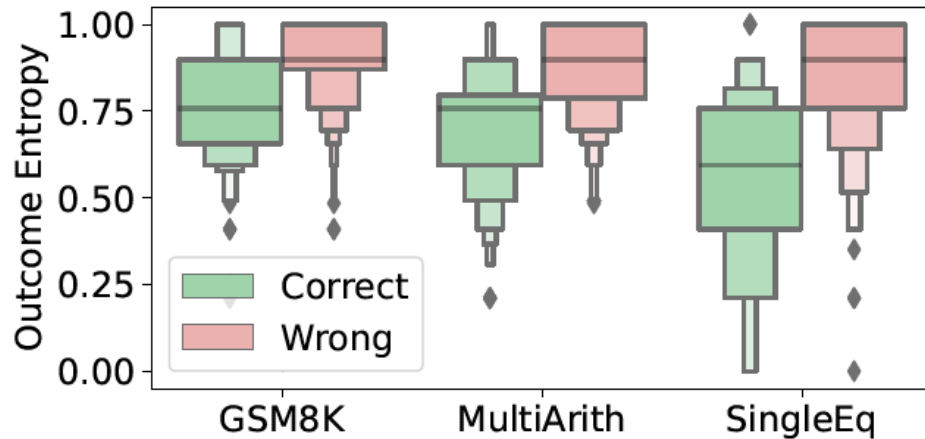
1. Self-consistency draws upon the insight that it approximates the amount of uncertainty (confidence) of the model for its prediction.

## Measurement (cont):

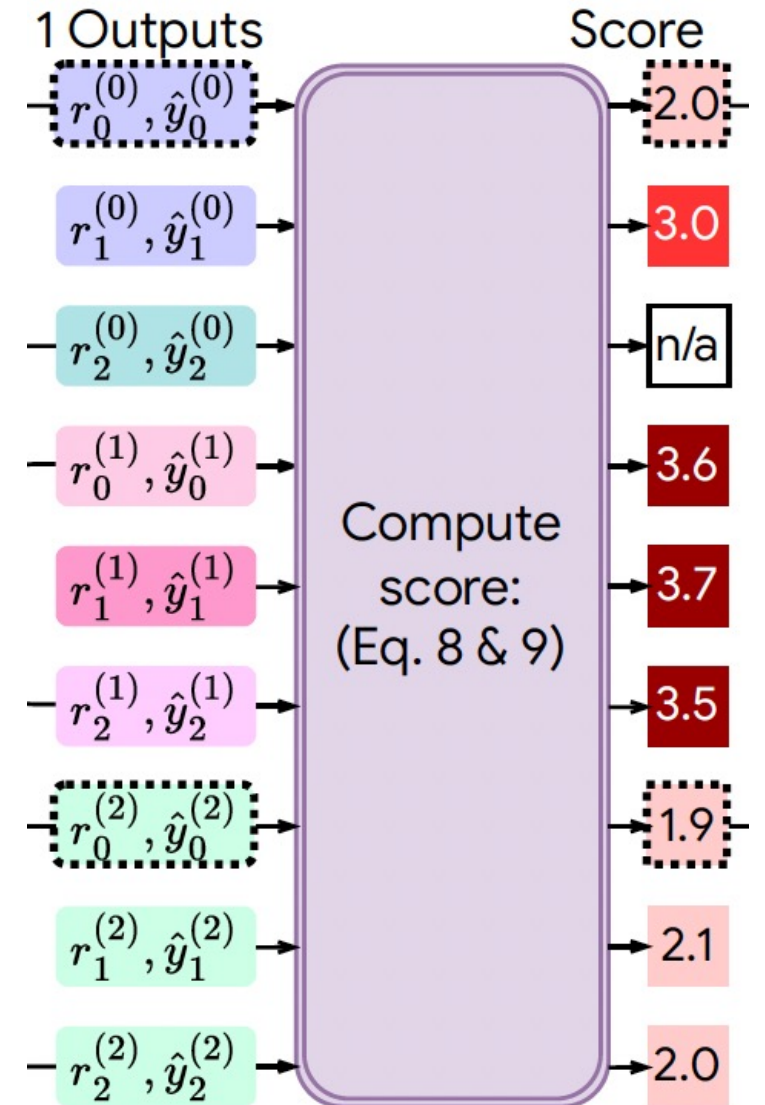
4. Compute the normalized entropy as:

$$\mathbb{H}(x^{(i)} | \{\hat{y}_j^{(i)}\}_{j=1}^m) = \frac{\sum_{\alpha=1}^n \tilde{p}(\hat{y}_\alpha^{(i)}) \log \tilde{p}(\hat{y}_\alpha^{(i)})}{\log m}, \quad (6)$$

where  $p$  is the empirical frequency of unique answer  $y_\alpha$ .



The **normalized entropy** is a good proxy over a number of different tasks where low entropy is positively correlated with correctness.





# Penalizing Repetitions

## Reasons:

1. “Repetitive demonstrations often lead to worse performance”.  
(Strong but spurious pattern)
2. Should capture semantic-level repetitions.

## Measurement:

1. Split demonstrations into phrases delimited by punctuation marks (“[.,?!]”).
2. Assuming with Q phrases, compute repetitiveness as:

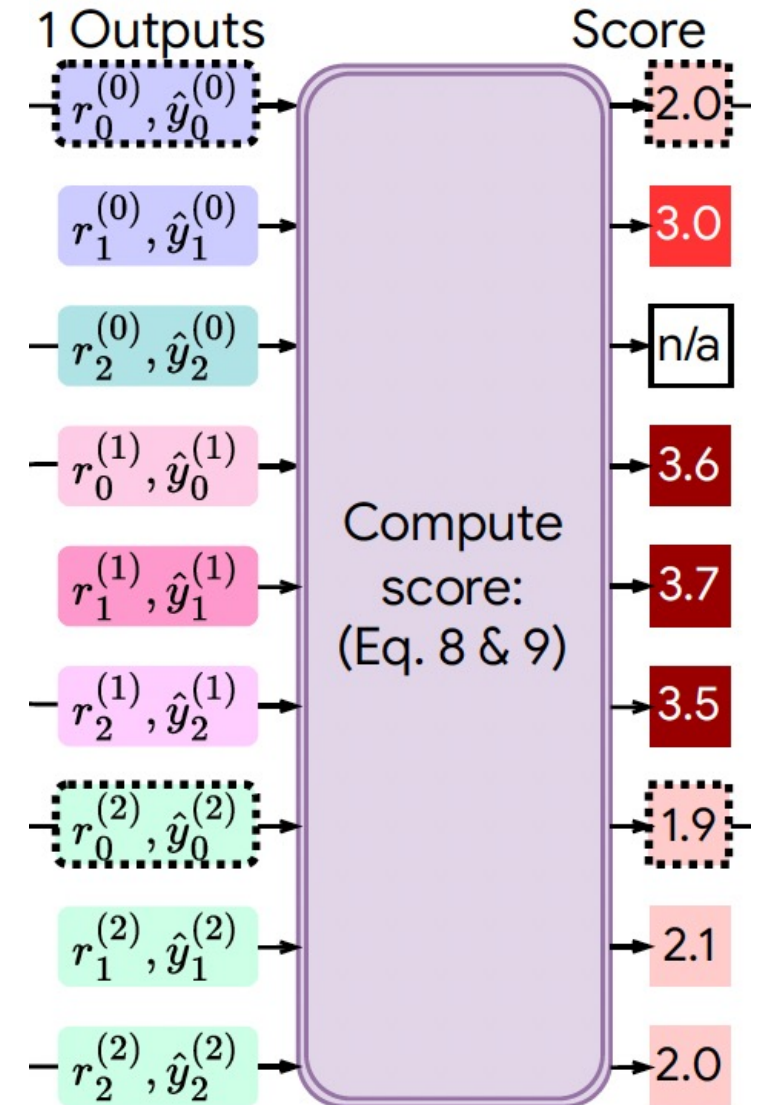
$$R_r(r_j^{(i)}) = \frac{2}{Q(Q-1)} \left( \sum_{a=1}^Q \sum_{b=a+1}^Q W_{ab} \right), \quad (7)$$

$$W_{ab} = S_c(\phi(q_a), \phi(q_b)),$$

where  $S_c(\cdot, \cdot)$  computes the cosine similarity and  $\phi(q_a)$  and  $\phi(q_b)$  denote the vector embedding of a-th and b-th phrases.

For now, the score is:

$$\mathcal{F}(p|x^{(i)}, r^{(i)}, \{\hat{y}_j^{(i)}\}_{j=1}^m) = \mathbb{H}(x^{(i)}) + \lambda R_r(r_j^{(i)}), \quad (8)$$



# Diversity

## Reasons:

1. To select a single in-context demonstration ( $K = 1$ ), we utilize the minimizer of the scoring function  $p^* = \arg \min_{p \in \mathcal{P}} \mathcal{F}(p)$
2. To select multiple demonstrations, we should penalize demonstrations that are similar to previous ones.

## Measurement:

1. Greedy forward selection with modified objective function:

$$\mathcal{G}_k(p) = \mathcal{F}(p) + \lambda R_q(p, \mathcal{S}_{k-1}) \quad \forall k \in [2, K], \quad (9)$$

$$R_q(p, \mathcal{S}_{k-1}) = \max \left( \{S_c(\phi(p), \phi(s_{k'}))\}_{k'=1}^{k-1} \right) \quad (10)$$

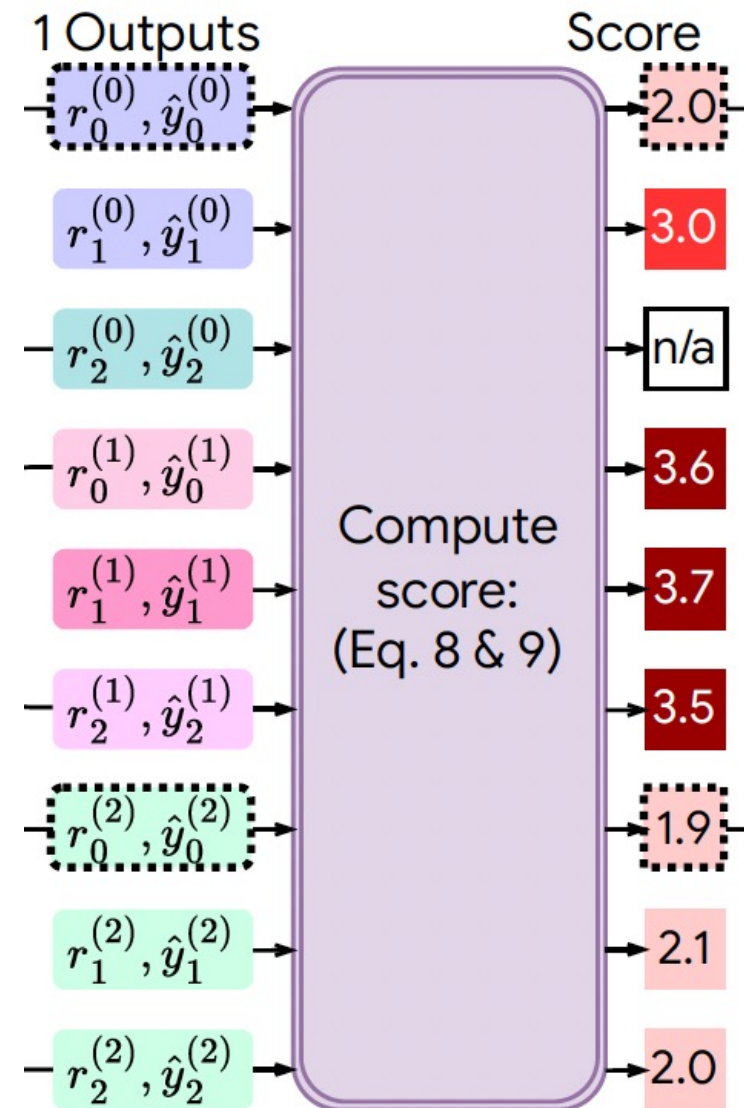
where  $\mathcal{S}_{k-1}$  is the partially built demonstration set  $\mathcal{S}$  with  $k - 1$  elements already selected.

---

### Algorithm 2 Building $\mathcal{S}$ for $K \geq 2$ .

---

- 1: Initialize  $\mathcal{S}$  with the **minimizer** of Eq. (8):  $\mathcal{S} \leftarrow \{p_0^* = \arg \min_{p \in \mathcal{P}} \mathcal{F}(p)\}$
  - 2: **for**  $k \in [2, K]$  **do**
  - 3: Find the minimizer of the modified objective (Eq. (9)):  $p_k^* = \arg \min_{p \in \mathcal{P}} \mathcal{G}_k(p)$ .
  - 4: Add  $p_k^*$  to  $\mathcal{S}$  and remove  $p_k^*$  from candidate pool  $\mathcal{P}$ .
  - 5: **end for**
- 



# Experiments

Model	PaLM-62B					PaLM-540B				
	Setting	0-shot		5-shot	Prev	Setting	0-shot		5-shot	Prev
Method	0-shot	Auto-	COSP	5-shot	8-shot	0-shot	Auto-	COSP	5-shot	8-shot
	CoT	CoT	(Ours)	CoT	CoT	CoT	CoT	(Ours)	CoT	CoT
# Paths	14	14	7+7	14		14	14	7+7	14	
MultiArith	67.2	9.4	<b>85.0</b>	<b>81.0</b>	-	95.2	<b>99.0</b>	<b>98.8</b>	96.0	99.3 <sup>b</sup>
AddSub	69.1	<b>73.2</b>	<b>78.9</b>	72.4	-	88.9	<b>89.1</b>	<b>89.9</b>	86.6	93.7 <sup>b</sup>
SingleEq	74.4	77.8	<b>78.7</b>	<b>79.8</b>	-	88.6	85.6	<b>90.4</b>	<b>89.2</b>	-
GSM-8K	20.9	9.2	<b>30.2</b>	<b>30.3</b>	27.4 <sup>a</sup>	68.5	71.4	<b>71.9</b>	64.3	74.4 <sup>b</sup>
CSQA	46.5	<b>68.2</b>	60.2	<b>66.8</b>	-	74.2	<b>79.4</b>	76.4	<b>80.7</b>	80.7 <sup>b</sup>
StrategyQA	57.2	59.4	<b>64.7</b>	<b>67.9</b>	-	66.0	<b>75.7</b>	75.2	<b>81.4</b>	81.6 <sup>b</sup>
(Average)	55.88	49.55	<b>66.28</b>	<b>66.37</b>	-	80.25	<b>83.37</b>	<b>83.77</b>	83.03	-

<sup>a</sup>Madaan and Yazdanbakhsh (2022). <sup>b</sup>Wang et al. (2022a): Significantly more (**40**) paths sampled.

Model	GPT-3 (code-davinci-001)				
	Setting	0-shot		5-shot	Prev
Method	0-shot	Auto-	COSP	5-shot	8-shot
	CoT	CoT	(Ours)	CoT	CoT
# Paths	14	14	7+7	14	
MultiArith	50.3	<b>78.5</b>	<b>80.7</b>	60.7	82.7 <sup>b</sup>
AddSub	43.5	31.9	<b>61.5</b>	<b>63.3</b>	67.8 <sup>b</sup>
SingleEq	48.8	58.1	<b>64.8</b>	<b>65.9</b>	-
GSM-8K	<b>10.2</b>	6.6	8.7	<b>16.7</b>	23.4 <sup>b</sup>
CSQA	29.2	50.9	<b>55.4</b>	<b>53.0</b>	54.9 <sup>b</sup>
StrategyQA	47.8	<b>55.8</b>	52.8	<b>55.4</b>	61.7 <sup>b</sup>
(Average)	38.32	46.96	<b>53.98</b>	<b>52.60</b>	-

<sup>b</sup>Wang et al. (2022a): Significantly more (**40**) paths sampled.

## Notes:

1. Unpretentious experimental design.
2. The outcome entropy is also a natural gauge of difficulty of questions to the LLM, as a higher entropy (thus a higher uncertainty) implies that the LLM may require additional demonstrations for this question.
3. Can further feature an adaptively allocated number of in-context demonstrations that is proportional to its zero-shot entropy in Stage 1, with higher-entropy questions given more demonstrations.

---

**Thank You!**