

# Rationale-Aware Answer Verification by Pairwise Self-Evaluation



Akira Kawabata\*

The Asahi Shimbun Company, Japan

Saku Sugawara

National Institute of Informatics, Japan

\* kawabata-a@asahi.com

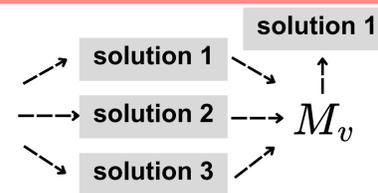
## TL;DR

- We propose a method to build verifiers that can verify **rationale validity** as well as answer correctness.
- Conventional training methods, which define correct-answer solutions as positive, fails to train rationale-aware verifiers, as **only 19% of model-generated correct-answer solutions have valid rationales**.
- We introduce iterative pairwise self-evaluation to identify solutions with high-quality rationales among correct-answer solutions.
- Our method improves verifiers' ability to detect valid rationale across multiple reasoning datasets.

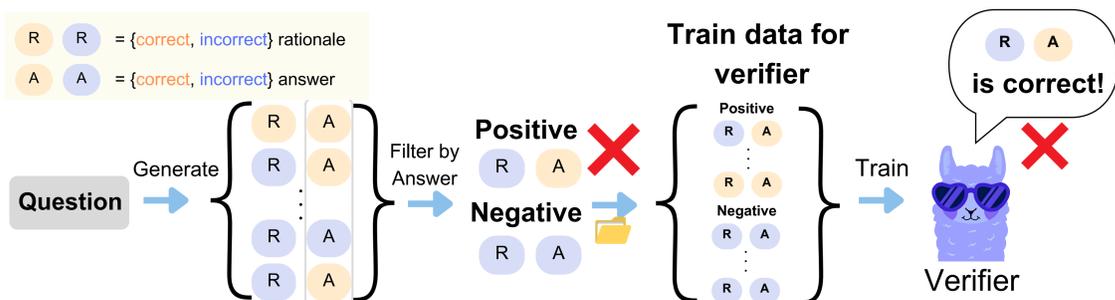
## Background & Motivation

### Q. What is verifier?

A. A model that evaluates the correctness  $M_g$  of model-generated solutions.



### Q. How is verifier (usually) trained?



A. Trained to distinguish positive and negative solutions. Conventionally, the labels are defined whether the solutions reach correct answers or not.

### Limitations

It is unclear whether the conventional approach can build rationale-aware verifiers.

### What We Did

We show that

- the conventional approach is insufficient for rationale-aware verifier due to the unfaithfulness in model-generated solutions.
- refining training data by iterative self-evaluation improves the verifiers' ability to distinguish valid rationales.

## Preliminary Experiments

### Settings

Model: Llama2-7B

Dataset: StrategyQA

**Rationale-quality annotation:** We use GPT-4 to annotate the validity of model-generated rationales.

### How often Does Model Generate Correct Answers with Valid Rationale?

#### Result

rationale ans **19%**

Only 19% of correct-answer solutions are judged as having valid rationales.

flawed rationale

### How Does Rationale Quality Affect Verifiers' Selection of Valid Rationales?

#### Method

#### Verifiers' Training Setting

Incrementally replace 10% of correct-answer solutions with valid solutions.

#### Eval Setting

For each Q, verifier selects from 5 solutions (1 valid, 2 correct-but-flawed, 2 incorrect)

$$\text{Rationale Accuracy} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbb{I}[M_v(q_i, \text{cands}_i) = \text{rationale}]$$

$$\text{Answer Accuracy} = \frac{1}{|D|} \sum_{i=1}^{|D|} \mathbb{I}[M_v(q_i, \text{cands}_i) = \text{ans}]$$

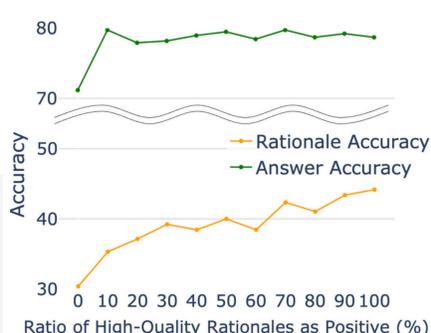
#### Result

As the proportion of valid solutions increases, Rationale Acc. improves while Answer Acc. remains largely stable.

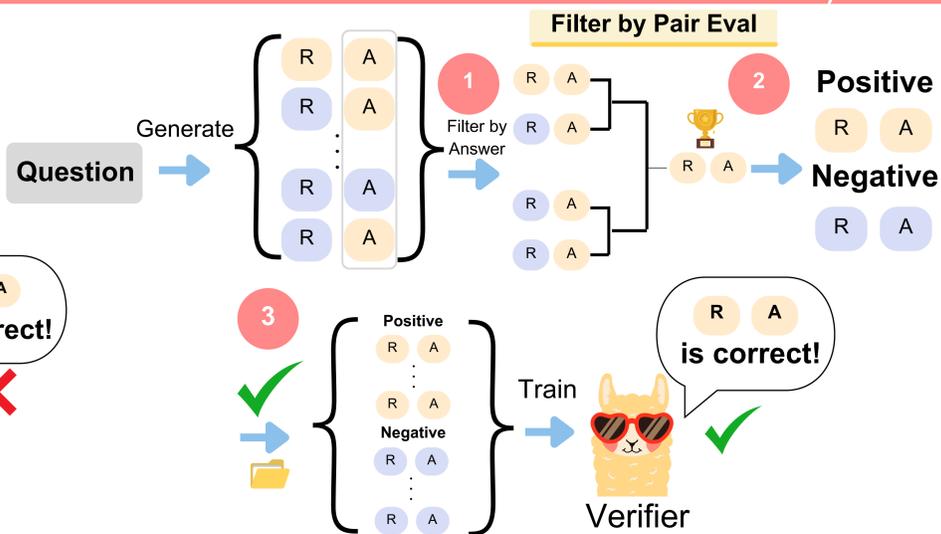
#### Implication

Large room for improving rationale accuracy in conventional verifier.

The key lies in how to identify high-quality model-generated rationales.



## Our Approach: REPS (Rationale Enhancement through Pairwise Selection)



- Selects  $N$  (e.g., 4) correct-answer solutions as candidates from model-generated solutions for each question.
- Applies tournament-style pairwise self-evaluation to these candidates, determining each round's winner through majority voting of  $S$  pairwise evaluations.
- Adds the final remaining solution to the positive samples for training the verifier.

## Experiments

### Settings

Model: Llama2-7B Dataset: ARC-Challenge, DROP, StrategyQA

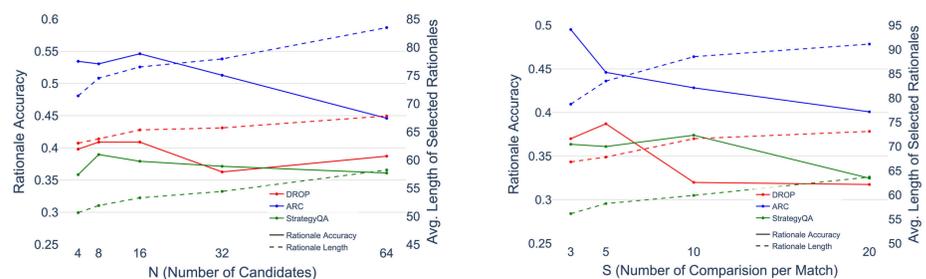
REPS Parameters:  $N=8, S=5$

Metrics: Rationale Accuracy, Task Performance (the accuracy of the verifier in selecting the answer from 5 generated solutions)

### Main Results

	ARC-Challenge		DROP		StrategyQA	
	Rationale Acc.	Task Perf.	Rationale Acc.	Task Perf.	Rationale Acc.	Task Perf.
Baseline	38.90	52.40	36.02	45.80	30.13	67.10
REPS	<b>53.05</b>	<b>54.75</b>	<b>40.90</b>	<b>46.90</b>	<b>38.96</b>	<b>67.25</b>

### Effects of Parameter Choice



- Unexpectedly, Rationale Accuracy decreases as  $N$  and  $S$  grow.
- This can be considered as the amplification of self-evaluation biases toward the superficial cue, i.e., the length of rationale.

### Which Verifier Can Select Better Rationales? Win Rate by GPT-4

	Result	Examples
DROP	213 (56.5%) REPS Wins, 164 (43.5%) Baseline Wins	Q. Is Disney associated with Los Angeles County? Disney has theme parks in Los Angeles County, California. Thus, Disney is associated with Los Angeles County. <b>REPS</b>
ARC	235 (60.4%) REPS Wins, 154 (39.6%) Baseline Wins	Disney is a large company that has many locations. Disneyland is located in Anaheim, which is 24 miles south of Los Angeles. Thus, Disney does not have a location in Los Angeles County. <b>Baseline</b>
StrategyQA	212 (58.1%) REPS Wins, 153 (41.9%) Baseline Wins	Q. A type of small mammal from the mountain regions of the western United States makes its home out of piles of rock. During summer months, the mammal places grasses and seeds in protected places in the rock piles. Which of the following is the most likely reason for this behavior? (correct answer: to store food that will be eaten over the winter months) This is the behavior of a prairie dog, which stores grass, seeds, and other food items in rock piles during summer months. Prairie dogs will consume the stored food during winter months. Thus, the prairie dog's behavior is to store food for use during the winter months. <b>REPS</b>

To protect and preserve its food, the mammal buries grasses and seeds in piles of rock. Grasses and seeds decay over time. Thus, the mammal buries its food to protect it from decay. **Baseline**