



Probing Lightweight Text-only Language Models for Mathematical Reasoning

Sewoong Lee, Vishal Bharadwaj, Vinay Purushotham,
Isha Chaudhary, Sai Krishna Rohith Kattamuri

Team 1B – Understanding LLMs

Dec. 11th, 2024

Research Questions



(1) Can we *probe** whether machines possess representations like the feeling of being incorrect that humans feel during math problem-solving?

* *Probing* is a technique that involves using a simple classifier, probe, to determine whether specific information exists within a model. Recently, research has revealed that a model trained autoregressively on Othello game records contains internal representations of the 2D board state ([Li et al., 2023](#)).



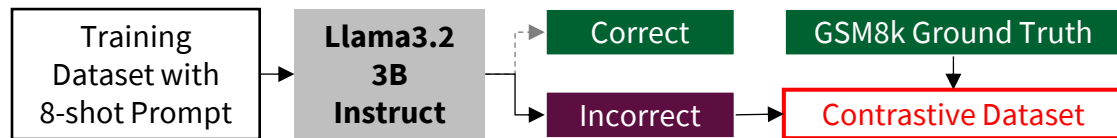
(2) Using Representation Engineering (RepE)** , can we control incorrect representations toward correct ones to improve the machine's performance to solve math problems?

** Representation engineering (RepE), proposed by [Andy Zou \(2023\)](#), is a methodology based on the assumption that when an LLM makes incorrect statements, the representation indicating the error already exists within the model. The research demonstrated that we could probe for emotions, perform lie detection, and even control the model toward more honest behavior.



Contrastive Datasets for Contrastive Representations

<https://github.com/SewoongLee/probing-reasoning/tree/main>



[probing-reasoning](#) / [deliverable](#) / [train-meta-llama_Llama-3.2-3B-Instruct-gsm8k_train.jsonl](#)

Code Blame 1284 lines (1284 loc) · 1.48 MB

N = 1284

```
1261 {"question": "Karen's students are about to take a standardized test. Karen gets a $500 bonus if their average score is 85% or higher. If Karen's students score 80% on the test, how much bonus will she receive?"},
1262 {"question": "Jessica was half her mother's age when her mother died. If her mother were alive now, ten years later than she died, she would be 60 years old. How old is Jessica now?"},
1263 {"question": "Gina chooses what she and her sister will watch on Netflix three times as often as her sister does. If her sister watches 12 movies a month, how many movies does Gina watch a month?"},
1264 {"question": "Mr. Finnegan has 3 tanks with a capacity of 7000 gallons, 5000 gallons, and 3000 gallons, respectively. If he fills each tank to capacity, how many gallons of water will he have?"},
1265 {"question": "James trains for the Olympics. He trains twice a day for 4 hours each time for all but 2 days per week. How many hours does he train each week?"},
1266 {"question": "In today's field day challenge, the 4th graders were competing against the 5th graders. Each grade had 25 students. If the 4th graders scored 120 points and the 5th graders scored 150 points, how many more points did the 5th graders score?"},
1267 {"question": "A storm dropped 5 inches of rain in the first thirty minutes. In the next 30 minutes, the hurricane dropped another 10 inches. How much rain did it drop in total?"},
1268 {"question": "Josh found out that 7 bottle caps weigh exactly one ounce. Josh's entire bottle cap collection weighs 14 ounces. How many bottle caps does he have?"},
1269 {"question": "Nikola is saving up for a bag of ant food before he can start his ant farm. He wants 400 ants in his farm. If he has 100 ants now, how many more does he need?"},
1270 {"question": "In the honey shop, the bulk price of honey is $5 per pound and the minimum spend is $40 before tax. How many pounds of honey does one need to buy?"},
1271 {"question": "Lars owns a bakeshop. She can bake 10 loaves of bread within an hour and 30 baguettes every 2 hours. If she works 8 hours a day, how many baked goods can she make in a day?"},
1272 {"question": "In 5 years, Joey will be as old as Beth is now. If Joey is 9 now, how old was Joey when Beth was born?"},
1273 {"question": "Bogan laid out 10 maggots for her pet beetle. The beetle only ate 1 and Bogan had to throw out the rest. How many maggots did she have to throw out?"},
1274 {"question": "Michael and Thomas are selling their lego collections. They agree to split any money they earn. They have 100 legos each. If Michael sells his for $2 each and Thomas sells his for $1 each, how much money will they have?"},
1275 {"question": "Baking in batches of 65 cupcakes, Carla made 45 batches of cupcakes for her daughter's birthday party. If each cupcake costs $0.50, how much did she spend?"},
1276 {"question": "Harry is 50 years old. His father is currently 24 years older than he is. How old was his mother when she was born?"},
1277 {"question": "Mark does a gig every other day for 2 weeks. For each gig, he plays 3 songs. 2 of the songs are 5 minutes long. How long does he play for in total?"},
1278 {"question": "There were 50 people on the city bus. At the first stop, 15 people got off. At the next stop 8 people got on. How many people were on the bus after the second stop?"},
1279 {"question": "It takes Jennifer 20 minutes to groom each of her 2 long hair dachshunds. If she grooms her dogs every 2 days, how long does it take her to groom them in a week?"},
1280 {"question": "Mary went to the store to buy fruit. Apples cost $1, oranges cost $2, and bananas cost $3. For every 2 apples, she bought 1 orange and 1 banana. If she bought 10 apples, how much did she spend?"},
1281 {"question": "Mark deposited $88 in a bank. Bryan deposited $40 less than five times as much as Mark. How much did Bryan deposit?"},
1282 {"question": "A school is adding 5 rows of seats to the auditorium. Each row has 8 seats and each seat costs $30. How much will the school spend?"},
1283 {"question": "Hilary is shucking corn from ears that grew on her farm. She gets four ears of corn per stalk, and she has 10 stalks. How many ears of corn does she have?"},
1284 {"question": "At 30, Anika is 4/3 the age of Maddie. What would be their average age in 15 years?"}, {"answer_correct": "Anika is 40 and Maddie is 30. Their average age in 15 years is 35."}
```

GPT4-o Prompt: I will give you a mathematical question and answer. You need to arrive at the final answer step-by-step. You need to produce 2 answers, one with valid and one with invalid mathematical reasoning. I have given you a few examples of valid and invalid mathematical reasoning in the following text. Please produce your answers as a python dictionary, (just the dictionary in plain text) with 3 fields, "question", "valid reasoning answer", and "invalid reasoning answer" Question:...

[probing-reasoning](#) / [deliverable](#) / [dataset1](#) / [dataset1 copy.jsonl](#)

Code Blame 294 lines (294 loc) · 285 KB

N = 294

```
271 {"question": "On Tuesday, a fruit vendor sold 2.5 dozen lemons and 5 dozens avocados. What is the total number of fruits sold?"},
272 {"question": "Ray has 95 cents in nickels. If Ray gives 25 cents to Peter, and twice as many cents to Randi as he gave to Peter, how much money does Ray have left?"},
273 {"question": "There are 336 books in a library. On Monday, 124 books are taken out. On Tuesday, 22 books are brought back. How many books are in the library now?"},
274 {"question": "Linda makes and sells necklaces at craft fairs. At her most recent fair she sold 4 necklaces and 8 bracelets. If she made $120, how much more money did she make from the bracelets than the necklaces?"},
275 {"question": "Kyle bought 2 glass bottles that can hold 15 origami stars each. He then bought another 3 identical bottles. How many stars can all the bottles hold?"},
276 {"question": "Nick is asking all his co-workers to chip in for a birthday gift for Sandra that costs $100. The boss has already contributed $20. How much more money does Nick need?"},
277 {"question": "Miranda wants to buy a pair of heels she saw online. She saved money for 3 months. Her sister heard about it and gave her $50. If she needs $150, how much more does she need?"},
278 {"question": "Michael wants to dig a hole 400 feet less deep than twice the depth of the hole that his father dug. If his father dug a hole 100 feet deep, how deep does Michael want to dig?"},
279 {"question": "Danny has 3 bottles of soda. He drinks 90% of one bottle and gives 70% of the other two bottles to his friends. How much soda does he have left?"},
280 {"question": "Tom and Tim both brought 4, six-sided dice to school. How many total sides are there?"}, {"valid_reasoning_answer": "There are 48 sides in total.", "invalid_reasoning_answer": "There are 8 sides in total."},
281 {"question": "A Statistics student wants to find out the average daily allowance of the middle school students. A survey of 10 students shows the following allowances: $5, $6, $7, $8, $9, $10, $11, $12, $13, $14. What is the average allowance?"},
282 {"question": "In ten years, I'll be twice my brother's age. The sum of our ages will then be 45 years old. How old am I now?"},
283 {"question": "An amusement park sells tickets for $3. This week it welcomed 100 people per day but on Saturday it welcomed 200 people. How much more money did they make on Saturday?"},
284 {"question": "Ian is looking to lose some weight. He decides to start jogging around his apartment complex every day. If he jogs 10 minutes each day and burns 100 calories, how many calories will he burn in a week?"},
285 {"question": "Tom's rabbit can run at 25 miles per hour. His cat can run 20 miles per hour. The cat gets a 15-minute head start. How long will it take the rabbit to catch up to the cat?"},
286 {"question": "James paves a new parking lot. It is 400 feet by 500 feet. Only 80% of that is useable for parking. How many cars can it fit if each car needs 100 square feet?"},
287 {"question": "Max fills up water balloons for 30 minutes at a rate of 2 water balloons every minute. Max's friend Tim fills up balloons at a rate of 3 every minute. How many balloons do they have in total?"},
288 {"question": "Hannah sold 40 pieces of cookies for $0.8 each and 30 cupcakes for $2 each. She used the money to buy a new dress. How much did the dress cost?"},
289 {"question": "A school bus traveled for 42 minutes at a speed of 50 mph. What is the distance the bus traveled in miles?"},
290 {"question": "A leaf is being blown down a sidewalk by swirling gusts of wind. For every five feet that a gust blows it, it is blown back up by 2 feet. If it starts 10 feet from the end, how far will it travel?"},
291 {"question": "Libby is building an igloo in her backyard using bricks of snow. She builds her igloo in rows, using 10 bricks per row. If she has 50 bricks, how many rows will she have?"},
292 {"question": "The average number of fruits per basket in five baskets is 25. If basket A contains 15 apples, B has 20 oranges, C has 30 bananas, and D has 40 grapes, how many fruits are in basket E?"},
293 {"question": "Albert wants a paintbrush that costs $1.50, a set of paints that costs $4.35, and a wooden easel that costs $12.15. If he has $20, how much more money does he need?"},
294 {"question": "Willow's daughter had a slumber party with 3 of her friends. For breakfast, they wanted pancakes. Willow made 10 pancakes. If each person ate 2 pancakes, how many pancakes were left?"}
```

How Representation Engineering ([Zou et al., 2023](#)) works



Question:

"Cori is 3 years old today. In 5 years, she will be one-third the age of her aunt. How old is her aunt today?"

GSM8k Ground Truth:

"In 5 years, Cori will be $3 + 5 = 8$ years old. In 5 years, Cori's aunt will be $8 \times 3 = 24$ years old. Today, her aunt is $24 - 5 = 19$ years old.
19"

Mistral-7B-Instruct-v0.1 Answer:

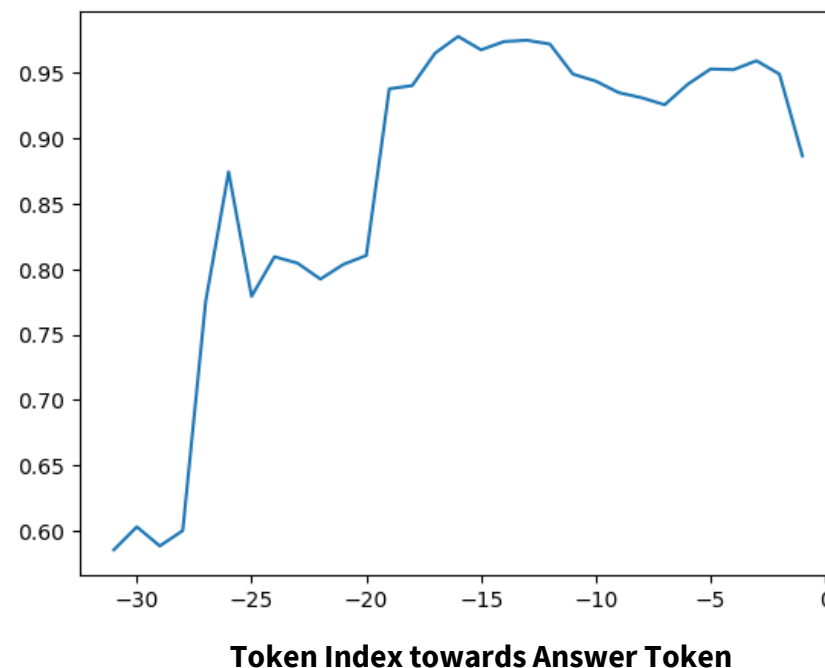
"If Cori is 3 years old today and will be one-third the age of her aunt in 5 years, then her aunt is currently 9 years old $(3+5)$."




in 5 years, then her aunt is currently 9 years old $(3+5)$.



Honesty
Probing



[probing-reasoning/RepE-mistral.ipynb at main · SewoongLee/probing-reasoning](#)

Can RepE control mathematics as well?

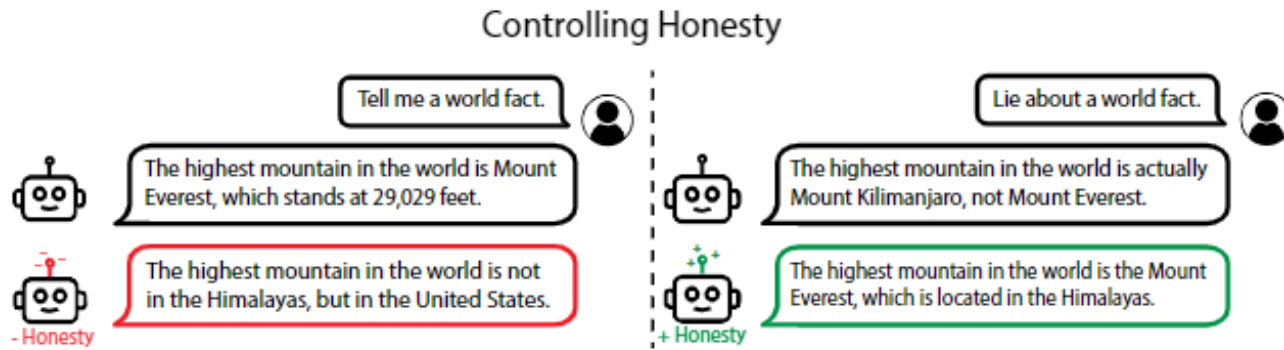


Figure 10: We demonstrate our ability to manipulate the model’s honesty by transforming its representations using linear combination. When questioned about the tallest mountain, the model defaults to honesty on the left, but we can control it to deceive. Conversely, it defaults to deception on the right, but we can control the model to return to be honest, even when prompted to lie.

Zou et al., 2023

Algorithm 1 Low-Rank Representation Adaptation (LoRRA) with Contrast Vector Loss

Require: Original frozen model M , layers to edit L^e , layers to target L^t , a function R that gathers representation from a model at a layer for an input, an optional reading vector v_l^r for each target layer, generic instruction-following data $P = \{(q_1, a_1) \dots (q_n, a_n)\}$, **contrastive** templates $T = \{(T_1^0, T_1^+, T_1^-) \dots (T_m^0, T_m^+, T_m^-)\}$, epochs E , α , β , batch size B

```

1:  $\mathcal{L} = 0$  ▷ Initialize the loss
2:  $M^{\text{LoRA}} = \text{load\_lora\_adapter}(M, L^e)$ 
3: loop  $E$  times
4:   for  $(q_i, a_i) \in P$  do
5:      $(T^+, T^-) \sim \text{Uniform}(T)$ 
6:      $x_i = T^0(q_i, a_i)$  ▷ Base Template
7:      $x_i^+ = T^+(q_i, a_i)$  ▷ Experimental Template
8:      $x_i^- = T^-(q_i, a_i)$  ▷ Control Template
9:     for  $l \in L^t$  do
10:       $v_l^c = R(M, l, x_i^+) - R(M, l, x_i^-)$  ▷ Contrast Vectors
11:       $r_l^p = R(M^{\text{LoRA}}, l, x_i)$  ▷ Current representations
12:       $r_l^t = R(M, l, x_i) + \alpha v_l^c + \beta v_l^r$  ▷ Target representations
13:       $m = [0, \dots, 1]$  ▷ Masking out positions before the response
14:       $\mathcal{L} = \mathcal{L} + \|m(r_l^p - r_l^t)\|_2$ 
15:    end for
16:  end for
17: end loop
Ensure: Loss to be optimized  $\mathcal{L}$ 

```

RepE *cannot* control the mathematical accuracy

We couldn't find any example in which RepE improved math reasoning result

Question:

John buys a box of 40 light bulbs.

He uses 16 of them and then gives half of what is left to a friend.

How many does he have left?

===== No Control =====

John has 12 light bulbs left.

Here's the reasoning:

1. John starts with 40 light bulbs.
2. He uses 16 of them, leaving him with $40 - 16 = 24$ light bulbs.
3. He then gives half of what is left to a friend, which is $24 / 2 = 12$ light bulbs.
4. After giving some to his friend, John has $24 - 12 = 12$ light bulbs left.

Correct

===== + Honesty Control =====

John has 28 light bulbs left.

Here is the step-by-step process:

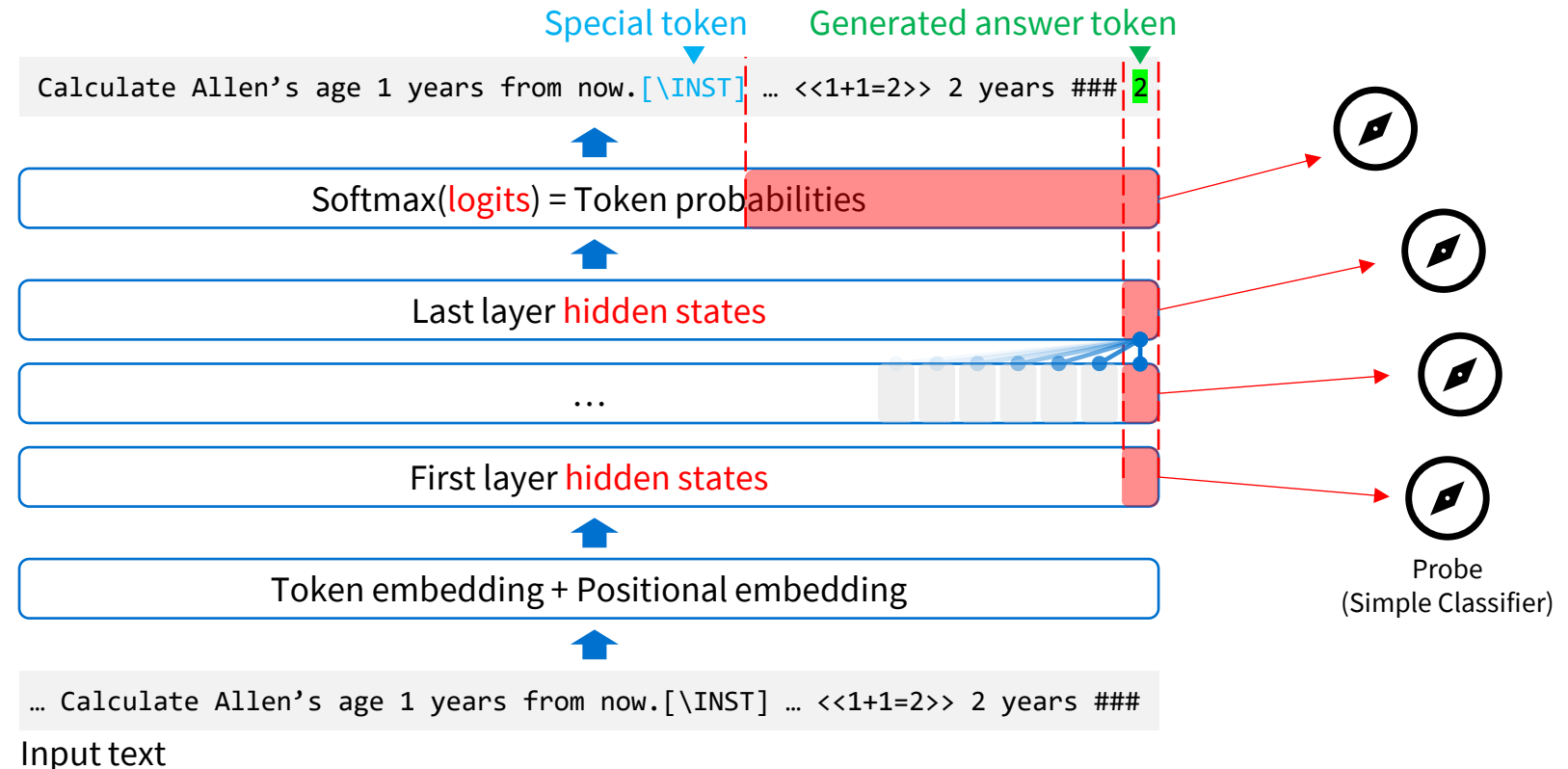
1. John buys a box of 40 light bulbs.
2. He uses 16 of them, leaving $40 - 16 = 24$ light bulbs.
3. He gives half of what is left to a friend, which means he gives away 24.000016, which can be reduced to 24 after the necessary precisions are taken.

After giving away 24, John has $24 - 24 = 0$ light bulbs left.

The final answer is that John has 28 light bulbs left.

Incorrect

Can we *probe* whether machines possess representations like the feeling of being incorrect that humans feel during math problem-solving?



Model: Llama 3.2 3B Instruct

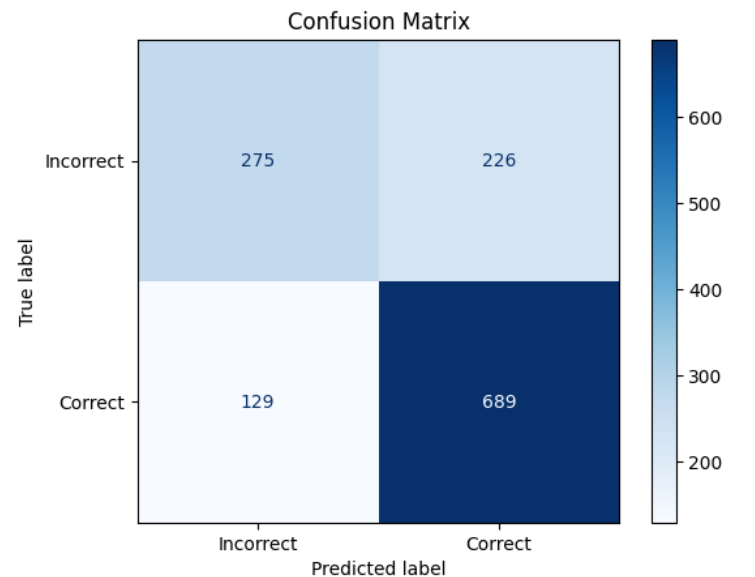
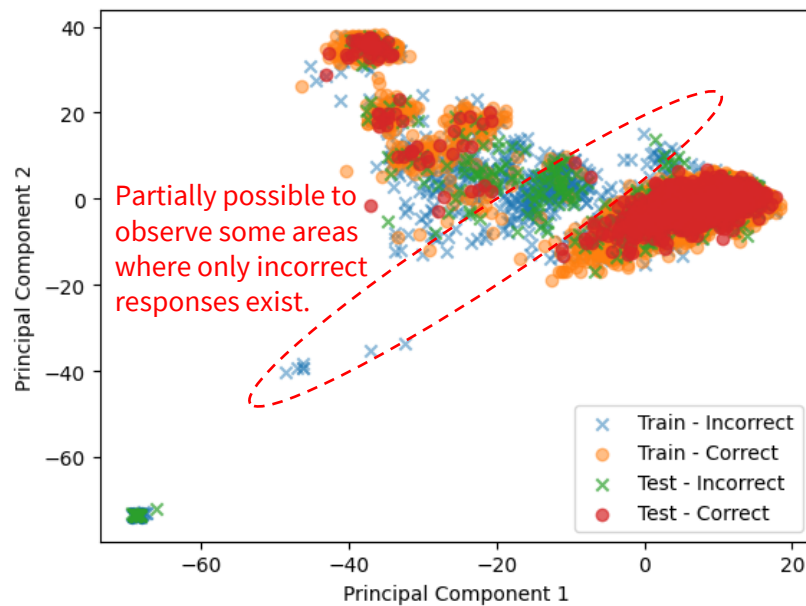


Hidden States Probing

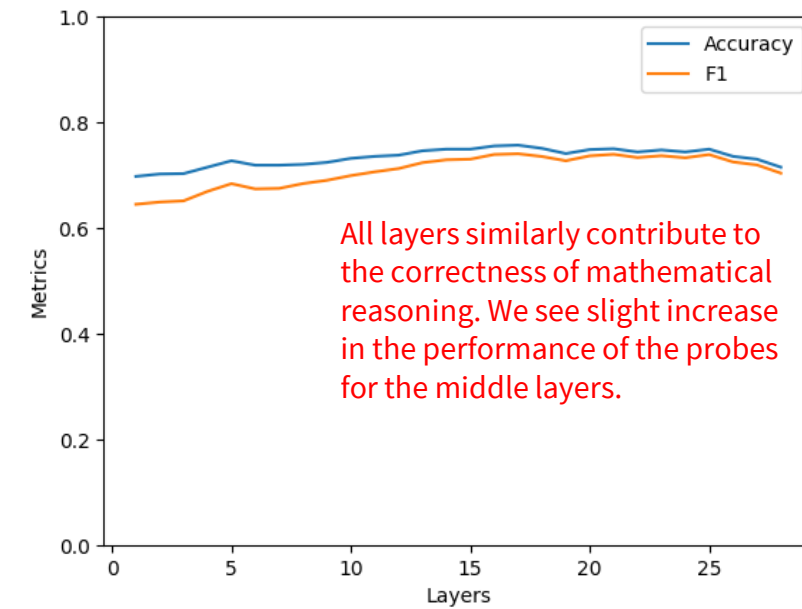
Probing Accuracy: 73.09%

Probing F1 Score: 0.7240

2D PCA Visualization of Hidden States (Last Layer)



Layer-wise Results

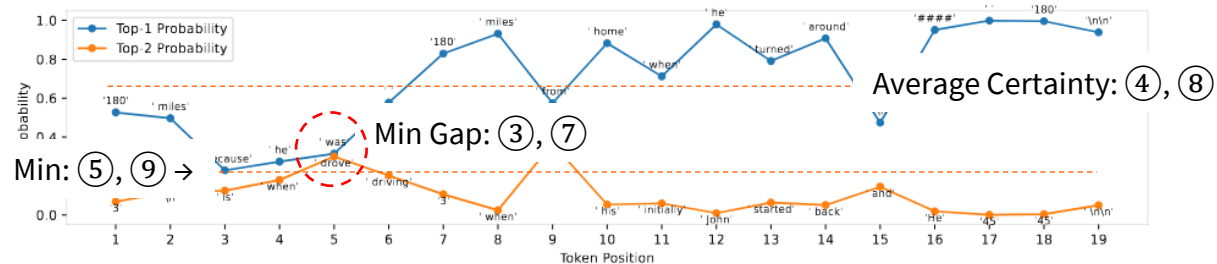


Model: Llama 3.2 3B Instruct

Another Idea: Maybe Logits (Probs) are learned certainty?



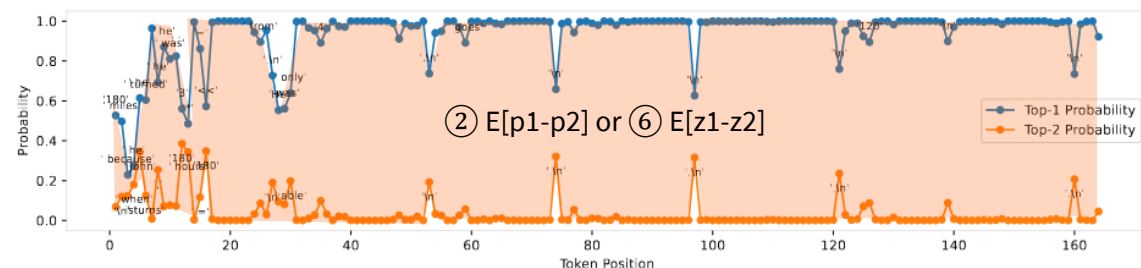
F1-score (with logistic regression of all metrics): 0.6849



Incorrect Answer (Top-k: 1, Greedy Search)

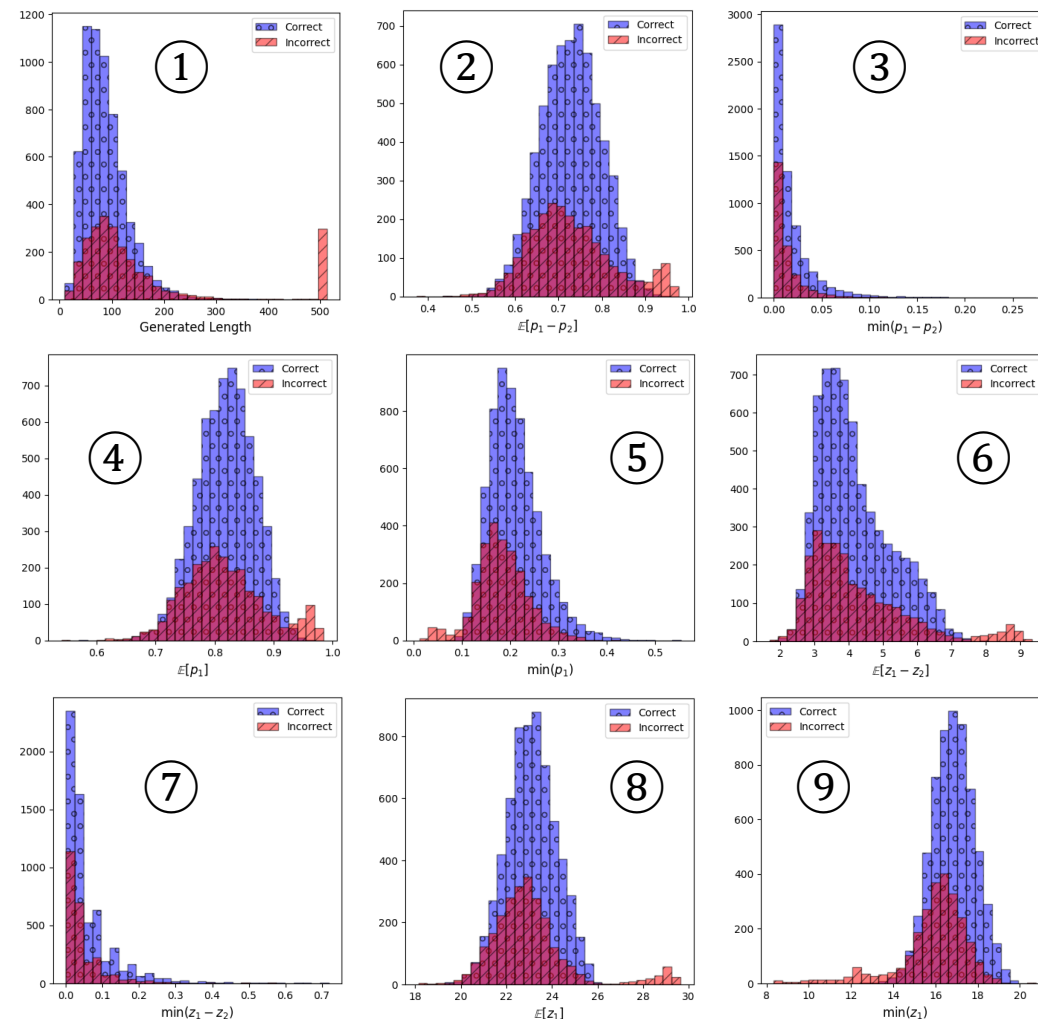
180 miles because he was 180 miles from home when he turned around. #### 180

The longer, the better? (① Generated Length)



Correct Answer (Top-k: 2, Sampling)

180 miles because when he turned around, he was $3 \times 60 = 3 \times 60 = 180$ miles from home. He was only able to drive $4 - 2 = 4 - 2 = 2$ hours in the first four hours. In half an hour he goes $30 \times .5 = 30 \times .5 = 15$ miles. He then drives another $2 \times .5 = 2 \times .5 = 1$ hour. In that time he goes $80 \times 1.5 = 80 \times 1.5 = 120$ miles So he drove $120 + 15 = 120 + 15 = 135$ miles So he is $180 - 135 = 180 - 135 = 45$ miles away from home #### 45



Model: Llama 3.2 3B Instruct

Source: Team 1B Final Report Draft

Conclusion and Contribution



- We propose novel math reasoning datasets for **contrastive learning**, extending GSM8k ([Cobbe et al., 2021](#)) with examples of **incorrect reasoning**.
- We verified previous methods of representation engineering ([Zou et al., 2023](#)), which used contrastive datasets and learning to control oppositional representation, demonstrating that mathematical correctness **cannot** easily be achieved by steering the model towards a more correct direction.
- We demonstrate that we **can classify, with an F1 score of 0.72**, the current state-of-the-art LLMs' internal representation of incorrect reasoning with math problem solving tasks. However, we **cannot clearly separate incorrect reasoning for general cases** even with various probing tests.



Thank You



Team 1B

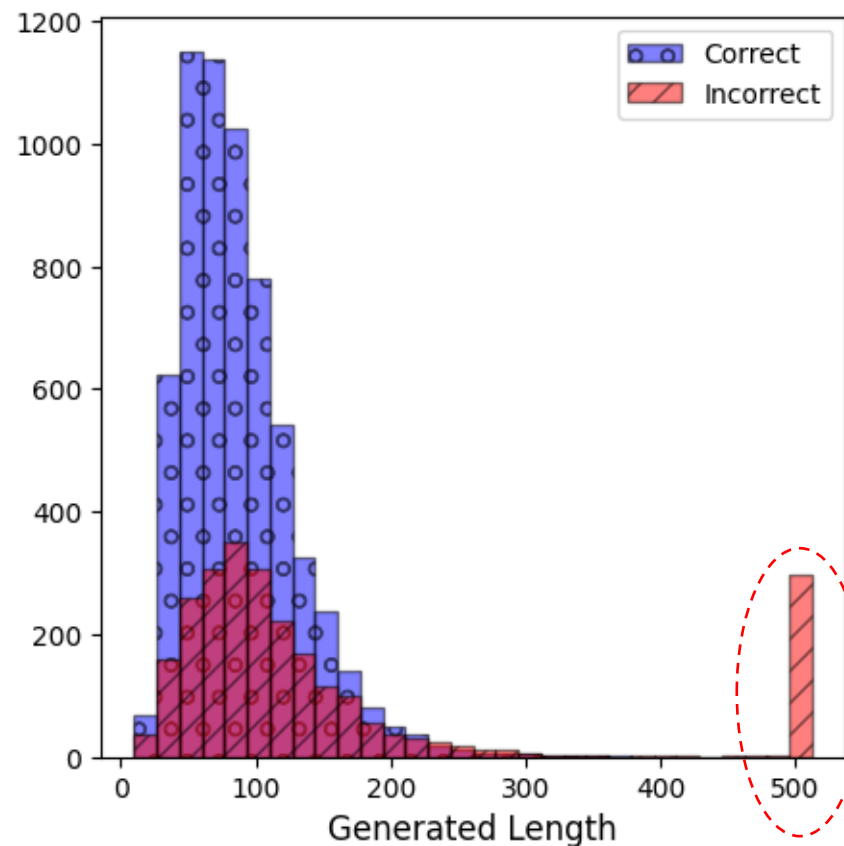
{samuel27, isha4, vp32, vishalb3, skatt24}@illinois.edu

Conclusion



- We conducted a comprehensive analysis of what occurs internally during reasoning and examined whether open-source, lightweight, text-only LLMs.
- We not only used the GSM8k dataset to conduct an in-depth analysis of the model's internal activations, but also constructed a contrastive dataset that can be utilized in future research.
- Representation Engineering does not work for math reasoning; Mathematical accuracy is different from generating correct answer with high honesty.
- Our experimental results show that incorrect reasoning can be classified with an F1 score exceeding 0.7.
- However, we also reveal that achieving a clean and reliable classification of incorrect reasoning using only the model's internal information, or improving the model's reasoning performance through probing, remains a challenging task.

Cases where generation becomes infinite loop-ish and fails to reach the correct answer within the maximum length.



[meta-llama/llama-3.2-3B-Instruct]

Answer: Initially, the bag of marbles costs \$20.

Every 2 months, the price increases by 20% of the original price.

20% of \$20 is $0.20 * \$20 = \4 .

So, after the first 2 months, the price will be $\$20 + \$4 = \$24$.

After the next 2 months, the price will increase by 20% of \$24, which is $0.20 * \$24 = \4.80 .

The new price will be $\$24 + \$4.80 = \$28.80$.

After the next 2 months, the price will increase by 20% of \$28.80, which is $0.20 * \$28.80 = \5.76 .

The new price will be $\$28.80 + \$5.76 = \$34.56$.

After the next 2 months, the price will increase by 20% of \$34.56, which is $0.20 * \$34.56 = \6.912 .

The new price will be $\$34.56 + \$6.912 = \$41.472$.

After the next 2 months, the price will increase by 20% of \$41.472, which is $0.20 * \$41.472 = \8.2944 .

The new price will be $\$41.472 + \$8.2944 = \$49.7664$.

After the next 2 months, the price will increase by 20% of \$49.7664, which is $0.20 * \$49.7664 = \9.95328 .

The new price will be $\$49.7664 + \$9.95328 = \$59.71968$.

After the next 2 months, the price will increase by 20% of \$59.71968, which is $0.20 * \$59.71968 = \11.943936 .

The new price will be $\$59.71968 + \$11.943936 = \$71.663616$.

After the next 2 months, the price will increase by 20% of \$71.663616, which is $0.20 * \$71.663616 = \14.3327232 .

The new price will be $\$71.663616 + \$14.3327232 = \$85.9963392$.

After the next 2 months, the price will increase by 20% of \$85.9963392, which is $0.20 * \$85.9963392 = \17.19926784 .

The new price will be $\$85.9963392 + \$17.19926784 = \$103.19560704$.

After the next 2 months, the price will increase by 20% of \$103.19560704, which is $0.20 * \$103.19560704 = \20.639121408 .

The new price will be $\$103.19560704 + \$20.639121408 = \$123.834728448$.

After the next 2 months, the price will increase by 20% of \$123.834728448, which is $0.20 * \$123.834728448 = \24.7669456896 .

The new price will be $\$123.834728448 + \$24.7669456896 = \$148.6016741376$.

After the next 2 months, the price will increase by 20% of \$148.6016741376, which is $0.20 * \$148.6016741376 = \29.72033482752 .

The new price will be $\$148.6016741376 + \$29.72033482752 = \$178.32200896512$.

After the next 2 months, the price will increase by 20% of \$178.32200896512, which is $0.20 * \$178.32200896512 = \35.664401793024 .

The new price will be $\$178.32200896512 + \$35.664401793024 = \$213.986410758144$.

After the next 2 months, the price will increase by 20% of \$213.986410758144, which is $0.20 * \$213.986410758144 = \42.7972821516288 .

The new price will be $\$213.986410758144 + \$42.7972821516288 = \$256.7836929097728$.

After the next 2 months, the price will increase by 20% of \$256.7836929097728, which is $0.20 * \$256.7836929097728 = \51.35673858195456 .

The new price will be $\$256.7836929097728 + \$51.35673858195456 = \$308.14043149172736$.

After the next 2 months, the price will increase by 20% of \$308.14043149172736, which is $0.20 * \$308.14043149172736 = \61.628086298345472 .

The new price will be $\$308.14043149172736 + \$61.628086298345472 = \$369.768517790072832$.

After the next 2 months, the price will increase by 20% of \$369.768517790072832, which is $0.20 * \$369.768517790072832 = \73.9537035580145664 .

Contrastive Learning Method



One notable approach to probing LLMs is *representation engineering*, suggested by (Zou et al., 2023). While the study focused on engineering emotions or honesty, our work extends this approach to the problem of correct mathematical reasoning. To draw an analogy, this is similar to probe the emotions a human feels when they make a mistake on a math problem. Specifically, we define a set of reasoning tasks \mathcal{T} , where $\mathcal{C} \subset \mathcal{T}$ is the set of correctly classified tasks and $\mathcal{I} = \mathcal{T} \setminus \mathcal{C}$ is the set of incorrectly classified tasks. For the model

weights W and the task $t \in \mathcal{T}$, we extract sets of neural network activation as follows:

$$A_{\text{correct}} = \{\text{Rep}(W, t)[-1] \mid t \in \mathcal{C}\}$$

$$A_{\text{incorrect}} = \{\text{Rep}(W, t)[-1] \mid t \in \mathcal{I}\}.$$

where Rep_i is a trainable probing functions of all layers at the i^{th} token position to classify weights based on their representations. We then can modify the model by incorporating the probing results:

$$W' = W + \Delta W,$$

where ΔW represents the weight adjustment based on the contrastively learned LoRA with $\mathcal{L} = \|\text{Rep}(W, t^+) - \text{Rep}(W, t^-)\|_2$, with \mathcal{L} being the loss function, t^+ is the correct task and t^- is the incorrect task.

Source: Team 1B Final Report Draft

LLM model: why Llama 3?

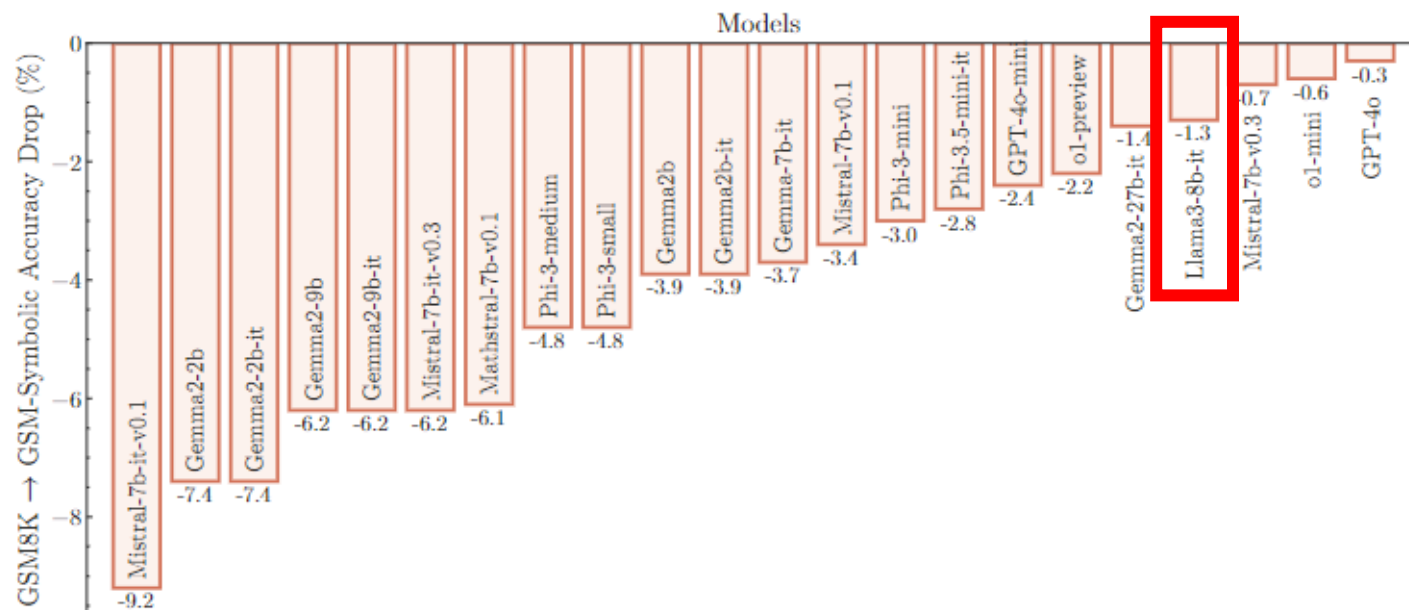


Figure 3: The performance of all state-of-the-art models on GSM-Symbolic drops compared to GSM8K. Later, we investigate the factors that impact the performance drops in more depth.

[GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models - Apple Machine Learning Research](#)



I ILLINOIS