

# MMReason: An Open-Ended Multi-Modal Multi-Step Reasoning Benchmark for MLLMs Toward AGI

## Supplementary Material

### A. Selected Benchmark details

In this section, we introduce several challenging benchmarks (*i.e.*, MMMU, MMMU-Pro, MMStar, M3CoT, MathVerse, and MathVision), which serve as the existing sources for questions that are reformulated into the open-ended format in MMReason.

**MMMU [53]:** MMMU is a large-scale benchmark with 11,550 questions spanning six major disciplines and 30 subjects, drawn from real college-level exams and materials. Covering a wide range of image types—including charts, maps, and diagrams—MMMU challenges multi-modal models to perform advanced reasoning and expert-level understanding. It serves as a rigorous testbed, as even leading models like GPT-4V and Gemini achieve only moderate accuracy, highlighting the benchmark’s difficulty.

**MMMU-Pro [54]:** MMMU-Pro is an advanced benchmark with 1,730 questions, extending MMMU to provide a stricter evaluation of multimodal reasoning. By filtering out text-only questions, expanding answer choices, and embedding queries within images, MMMU-Pro requires models to simultaneously process visual and textual information. This setting yields lower model performance than MMMU, offering a challenging and realistic testbed for multimodal understanding and reasoning.

**MMStar [3]:** MMStar is a vision-centric multimodal benchmark consisting of 1,500 challenging samples. It is designed to evaluate six core capabilities—coarse perception, fine-grained perception, instance reasoning, logical reasoning, science & technology, and mathematics—across 18 specific dimensions, offering a balanced and comprehensive assessment of MLLMs’ multimodal reasoning abilities.

**M3CoT [5]:** M3CoT features 11,459 problems across multiple domains, targeting multi-modal chain-of-thought reasoning that integrates both textual and visual information in multi-step processes. By overcoming the limitations of prior benchmarks, M3CoT offers a diverse and challenging platform for evaluating MLLMs and reveals the ongoing gap between VLLM and human performance in complex reasoning tasks.

**MathVerse [58]:** MathVerse is a visual math benchmark comprising 2,612 multi-subject problems with diagrams, drawn from public datasets. Each problem is reformulated by annotators into six types: Text Dominant, Text Lite, Text Only, Vision Intensive, Vision Dominant, and Vision Only. For this study, we assess StepGRPO on five vision-related versions, excluding Text Only.

**MathVision [37]:** MathVision features 3,040 rigorously

selected math questions enriched with visual elements, all derived from authentic mathematics competitions. Spanning 16 mathematical fields and five difficulty tiers, it provides a varied and thorough basis for evaluating the mathematical reasoning skills of MLLMs.

### B. More Result on MMReason

We report results on the MMReason test set in Table 6, which show a similar trend to those observed on the test-mini set in Table 3.

### C. Additional qualitative analysis of reasoning

We provide additional qualitative results of the powerful GPT-4o’s reasoning responses, including both correct and incorrect ones. For incorrect responses, we conduct a detailed analysis of their underlying causes, as illustrated from Figure 6 to Figure 24.

Models	Size	Final	Discipline-wise final reasoning accuracy					
		Acc.	Math	Business	Science	Engineering	Social	Health
Closed-Source Models								
Gemini-2.5-pro-preview-03-25 [6]	-	48.2	52.4	55.4	36.2	32.3	45.4	65.8
Claude-sonnet-4-20250514-thinking [32]	-	38.6	44.2	44.9	23.2	20.5	40.2	58.2
GPT-4.1 [27]	-	30.9	28.6	47.3	26.2	15.2	39.2	58.2
GPT-4o-1120 [17]	-	24.5	22.4	42.2	15.8	11.5	35.1	51.9
Open-Source Models								
Deepseek-VL2 [43]	4.1B <sup>†</sup>	11.4	11.1	14.1	7.3	5.6	28.9	31.6
LLaVA-OneVision [20]	7B	11.5	12.6	14.4	5.7	3.1	24.7	29.1
Mulberry [46]	7B	12.9	13.1	17.8	7.8	5.6	23.7	26.6
Qwen-2.5-VL [38]	7B	19.1	18.9	26.4	14.2	9.9	25.8	36.7
MiniCPM-V-2.6 [49]	8B	11.7	12.5	13.3	8.1	4.0	27.8	22.8
InternVL-2.5-MPO [38]	8B	15.4	17.2	17.1	9.7	6.2	27.8	27.8
LLaMA-3.2-Vision [8]	11B	9.7	8.9	12.6	7.6	3.1	24.7	29.1
InternVL-2.5-MPO [38]	78B	27.3	26.7	40.1	18.9	11.8	30.9	51.9
Qwen-2.5-VL [2]	32B	25.2	27.1	35.1	16.1	9/0	29.9	44.3
Qwen-2.5-VL [2]	72B	26.9	26.9	43.2	18.1	9.3	29.9	48.1

Table 6. **MMReason result** on test set. We provide the final-answer reasoning accuracy for 13 mainstream MLLMs, as well as the final answer reasoning accuracy across different disciplines. <sup>†</sup> denotes activated parameters of MoE model.

### Correct example

#### Question

The following information in the image was taken from the records of Katana Limited. New equipment was purchased during the year at a cost of R7 000 000. Old equipment was sold at carrying value. Calculate the returns on non-current assets sold.

	28 FEBRUARY 2021	28 FEBRUARY 2020
Non-current assets at cost	18 000 000	12 000 000
Accumulated depreciation	1 440 000	1 200 000
Depreciation for the year	620 000	430 000

#### GPT-4o

To determine the returns on non-current assets sold, we need to calculate the carrying value (cost minus accumulated depreciation) and analyze the sale proceeds. The question states that the old equipment was sold at **carrying value**, meaning the sale price equals the carrying value.

Information provided:

- Non-current assets at cost (opening balance, 28 Feb 2020): **R12,000,000**
- Non-current assets at cost (closing balance, 28 Feb 2021): **R18,000,000**
- Accumulated depreciation (opening balance, 28 Feb 2020): **R1,200,000**
- Accumulated depreciation (closing balance, 28 Feb 2021): **R1,440,000**
- Depreciation for the year: **R620,000**
- New equipment purchased: **R7,000,000**
- Old equipment sold at **carrying value**.

Step 1: Calculate the cost of old equipment sold.

The movement in the **Non-current assets at cost** account reflects:

Closing Balance = Opening Balance + Additions – Disposals

Substituting values:

$18,000,000 = 12,000,000 + 7,000,000 - \text{Disposals at Cost}$

$\text{Disposals at Cost} = 12,000,000 + 7,000,000 - 18,000,000 = 1,000,000$

So, the **cost of old equipment sold = R1,000,000**.

Step 2: Calculate the accumulated depreciation on old equipment sold.

The movement in the **Accumulated Depreciation** account reflects:

Closing Balance = Opening Balance + Depreciation for the Year - Accumulated Depreciation on Disposals

Substituting values:

$1,440,000 = 1,200,000 + 620,000 - \text{Accumulated Depreciation on Disposals}$

$\text{Accumulated Depreciation on Disposals} = 1,200,000 + 620,000 - 1,440,000 = 380,000$

So, the accumulated depreciation on old equipment sold = R380,000.

Step 3: Calculate the carrying value of the old equipment sold.

Carrying value = Cost - Accumulated depreciation:

Carrying value =  $1,000,000 - 380,000 = 620,000$

Since the old equipment was sold at **carrying value**, the sale proceeds were **R620,000**, which matches the carrying value.

Final Answer:

The **returns on non-current assets sold** refer to proceeds received from the disposal of assets. Since the equipment was sold at **carrying value**, the returns = **R620,000**.

Figure 6. Qualitative analysis.

### Correct example

#### Question

In a case-control study, the results were shown in the table below. The OR was...

	Case group	Control group
A history of exposure	400	100
No history of exposure	100	400
Total	500	500

#### GPT-4o

To calculate the **odds ratio (OR)** for the given case-control study, we use the formula:

$$\text{OR} = \frac{(a \cdot d)}{(b \cdot c)}$$

Where:

- a = Number of cases with exposure = 400
- b = Number of controls with exposure = 100
- c = Number of cases without exposure = 100
- d = Number of controls without exposure = 400

Substitute the values into the formula:

$$\text{OR} = \frac{(400 \cdot 400)}{(100 \cdot 100)} = \frac{160,000}{10,000} = 16$$

The odds ratio (OR) is 16.

**Final Answer:**

This means that the odds of exposure are 16 times higher in the case group compared to the control group.

Figure 7. Qualitative analysis.

### Wrong example

#### Question

The following balances in the image were taken from the trial balance on 31 July 2021, the last day of the financial year. The mortgage loan was increased with R100 000 on 30 April 2021. This transaction was properly recorded in the financial records. Calculate the total interest amount for the year ending 31 July 2021.

Loan: Nedbank 10% p.a.	B8		350 000
Fixed deposit: Nedbank 8.50% p.a.	B11	120 000	

#### GPT-4o

To calculate the total interest for the year ending 31 July 2021, we need to consider the interest on both the mortgage loan and the fixed deposit.

##### Step 1: Interest on the loan

1. **Loan balance before April 30, 2021**: The loan balance as of 1 August 2020 was **R350,000**, but it increased by **R100,000** on **30 April 2021**. This means the revised loan balance became **R450,000** from that date onward. We will calculate the interest separately for the periods before and after the increase.

**[Question Misunderstanding]** The revised loan balance is R350,000. It was R250,000 previously.

Interest on loan (split into two periods):

- **Period 1** (1 August 2020 to 30 April 2021):

Loan balance = R350,000

Interest rate = 10% per annum

Number of months = 9 months

Interest =  $(350,000 \times 0.10 \times \frac{9}{12})$

Interest = R26,250

- **Period 2** (1 May 2021 to 31 July 2021):

Loan balance = R450,000

Interest rate = 10% per annum

Number of months = 3 months

Interest =  $(450,000 \times 0.10 \times \frac{3}{12})$

Interest = R11,250

- **Total loan interest**:

$R26,250 + R11,250 = R37,500$

**[Question Misunderstanding]** Fixed deposit is not a mortgage loan. This should not be included.

##### Step 2: Interest on the fixed deposit

The fixed deposit balance is **R120,000**.

Interest rate = 8.5% per annum.

The deposit amount is fixed throughout the year, so we calculate interest for the full year.

Interest =  $(120,000 \times 0.085)$

Interest = **R10,200**

##### Step 3: Total interest calculation

Now, sum up the loan interest and the fixed deposit interest to find the total interest for the year ending 31 July 2021.

Total interest = Loan interest + Fixed deposit interest

Total interest =  $(R37,500 + R10,200)$

Total interest = **R47,700**

##### Final Answer:

The total interest amount for the year ending 31 July 2021 is **R47,700**.

Figure 8. Qualitative analysis.

### Wrong example

#### Question

The following information in the image was taken from the loan account statement received for the year ending 28 February 2020.

What was the balance of the loan account on 1 March 2019?

Balance on 1 March 2019	?
Fixed monthly payment of R9 600 (including interest)	?
Capitalised interest for the year	R56 700
Balance on 28 February 2020	R482 000

#### GPT-4o

To calculate the balance of the loan account on **1 March 2019**, we need to account for the **fixed monthly payments** and the **capitalized interest for the year**. Here's how to compute it step by step:

### Known values from the question:

- **Fixed monthly payment**: R9,600 (includes both interest and loan repayment)
- **Capitalized interest for the year**: R56,700
- **Balance on 28 February 2020**: R482,000

### Formula:

The relationship between the balances on the loan account is:

Opening Balance on 1 March 2019 = Closing Balance on 28 February 2020 + Loan Repayments during the year - Capitalized Interest for the year

**Step 1: Calculate total loan repayments during the year.**

Since the fixed monthly payment is **R9,600** and there are **12 months** in the year, the total payments during the year are:

$$\text{Total Payments} = R9,600 \times 12 = R115,200$$

**Step 2: Separate the loan repayment portion from the total payments.**

Out of the total payments, some portion covers interest, while the remainder goes to reducing the loan principal. The **capitalized interest** for the year is **R56,700**, so the portion of the payments used to reduce the loan balance is:

$$\text{Loan Repayment Portion} = \text{Total Payments} - \text{Capitalized Interest}$$

$$\text{Loan Repayment Portion} = R115,200 - R56,700 = R58,500$$

**Step 3: Compute the opening balance.**

Using the formula above:

$$\text{Opening Balance on 1 March 2019} = \text{Closing Balance on 28 February 2020} + \text{Loan Repayment Portion} - \text{Capitalized Interest}$$

Substitute the known values:

$$\text{Opening Balance on 1 March 2019} = R482,000 + R58,500 - R56,700$$

$$\text{Opening Balance on 1 March 2019} = R483,800$$

**[Conceptual Confusion]** Capitalized interests was already deducted in Step 2. It should not be deducted again here.

**Final Answer:**

The balance of the loan account on **1 March 2019** was **R483,800**.  $R482,000 + R58,500 = R540,500$

Figure 9. Qualitative analysis.

### Wrong example

#### Question

Suppose that your demand schedule for DVDs is as follows: <image 1> what is your price elasticity of demand as the price of DVDs increases from \$8 to \$10 if your income is \$12,000?

Price	Quantity Demanded(income = \$10,000)	Quantity Demanded(income = \$12,000)
\$ 8	40 DVDs	50 DVDs
10	32	45
12	24	30
14	16	20
16	8	12

#### GPT-4o

To calculate the \*\*price elasticity of demand (PED)\*\*, we use the following formula:

$$PED = \frac{\% \text{Change in Quantity Demanded}}{\% \text{Change in Price}}$$

#### Step 1: Calculate the percentage change in quantity demanded

From the table, when the price increases from \$8 to \$10 and income is \$12,000, the quantity demanded decreases from \*\*50 DVDs to 45 DVDs\*\*.

$$\begin{aligned} \text{Change in Quantity Demanded} &= \frac{\text{New Quantity} - \text{Old Quantity}}{\text{Old Quantity}} \times 100 \\ \text{Change in Quantity Demanded} &= \frac{45 - 50}{50} \times 100 = \frac{-5}{50} \times 100 = -10\% \end{aligned}$$

#### Step 2: Calculate the percentage change in price

The price increases from \$8 to \$10.

$$\begin{aligned} \text{Change in Price} &= \frac{\text{New Price} - \text{Old Price}}{\text{Old Price}} \times 100 \\ \text{Change in Price} &= \frac{10 - 8}{8} \times 100 = \frac{2}{8} \times 100 = 25\% \end{aligned}$$

#### Step 3: Calculate the price elasticity of demand

Now, substitute the values into the PED formula:

$$PED = \frac{-10\%}{25\%} = -0.4$$

The absolute value of PED is:  $|PED| = 0.4$

#### Final Answer:

The price elasticity of demand as the price of DVDs increases from \$8 to \$10 (when income is \$12,000) is \*\*0.4\*\*. This indicates that the demand is \*\*inelastic\*\*.

**[Lack of Knowledge]** The wrong equations are used for both. Below are the correct ones.

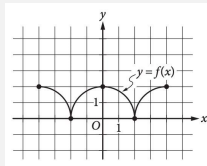
$$\% \Delta Q = \frac{Q_2 - Q_1}{\frac{Q_1 + Q_2}{2}} \times 100 \quad \% \Delta P = \frac{P_2 - P_1}{\frac{P_1 + P_2}{2}} \times 100$$

Figure 10. Qualitative analysis.

### Wrong example

#### Question

The figure above shows the complete graph of the function  $f$  in the  $x$ - $y$ -plane. The function  $g$  (not shown) is defined by  $g(x) = f(x) + 6$ . What is the maximum value of the function  $g$ ?



#### GPT-4o

To determine the maximum value of the function ( $g(x) = f(x) + 6$ ), we need to analyze the graph of ( $f(x)$ ) and find its maximum value, then add 6 to it.

**Image Information Mistake: The graph's highest point has a value of 2, yet its y-value is misinterpreted as 1, probably because of the marked cell height of 1.**

Step 1: Identify the maximum value of ( $f(x)$ )

From the graph of ( $f(x)$ ), the highest point on the curve is at ( $y = 1$ ). Therefore, the maximum value of ( $f(x)$ ) is:

$$\text{Maximum of } f(x) = 1$$

Step 2: Calculate the maximum value of ( $g(x)$ )

The function ( $g(x)$ ) is defined as ( $g(x) = f(x) + 6$ ). To find the maximum value of ( $g(x)$ ), we add 6 to the maximum value of ( $f(x)$ ):

$$\text{Maximum of } g(x) = \text{Maximum of } f(x) + 6 = 1 + 6 = 7$$

Final Answer:

The maximum value of the function ( $g(x)$ ) is 7.

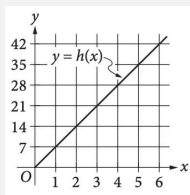
Figure 11. Qualitative analysis.



### Correct example

#### Question

The line in the  $x$   $y$ -plane above represents the relationship between the height  $h(x)$ , in feet, and the base diameter  $x$ , in feet, for cylindrical Doric columns in ancient Greek architecture. How much greater is the height of a Doric column that has a base diameter of 5 feet than the height of a Doric column that has a base diameter of 2 feet?



#### GPT-4o

To solve the problem, we need to determine the height of the Doric column for base diameters of 5 feet and 2 feet, and then find the difference between these heights.

##### Step 1: Understand the relationship

The graph shows a linear relationship between the base diameter ( $x$ ) and the height ( $h(x)$ ). The equation of the line is given as ( $y = h(x)$ ). From the graph, we can determine the slope and the  $y$ -intercept of the line.

##### Step 2: Find the slope of the line

The slope of the line is calculated as:

$$\text{slope} = \frac{\text{change in } y}{\text{change in } x}$$

From the graph, when ( $x = 1$ ), ( $y = 7$ ), and when ( $x = 2$ ), ( $y = 14$ ). Thus:

$$\text{slope} = \frac{14 - 7}{2 - 1} = 7$$

##### Step 3: Write the equation of the line

The line passes through the origin  $((0, 0))$ , so the  $y$ -intercept is 0. The equation of the line is:

$$h(x) = 7x$$

##### Step 4: Calculate the heights

- For a base diameter of ( $x = 5$ ):

$$h(5) = 7 \cdot 5 = 35, \text{ feet}$$

- For a base diameter of ( $x = 2$ ):

$$h(2) = 7 \cdot 2 = 14, \text{ feet}$$

##### Step 5: Find the difference in heights

The difference in heights is:

$$h(5) - h(2) = 35 - 14 = 21, \text{ feet}$$

##### Final Answer:

The height of a Doric column with a base diameter of 5 feet is 21 feet greater than the height of a Doric column with a base diameter of 2 feet.

Figure 12. Qualitative analysis.

### Correct example

#### Question

In 2018, which nation boasts the highest GDP per capita among the Country 3 and Country 2 countries?

Year	Country	GDP(\$)	Population (hundred millions)	Import Value (\$)
2017	Country 2	17.61	6.51	1.58
2018	Country 2	2.66	14.92	3.91
2018	Country 3	41.34	15.92	2.28
2021	Country 4	28.87	15.49	1.24
2021	Country 5	4.65	15.74	4.98
2022	Country 1	26.5	10.27	3.76
2022	Country 4	38.05	16.53	3.52
2023	Country 4	14.91	4.83	4.57

#### GPT-4o

To determine which country has the highest GDP per capita in 2018 between Country 3 and Country 2, we calculate GDP per capita using the formula:

$$GDP \text{ per capita} = \frac{GDP}{Population}$$

For Country 2 in 2018:

- GDP = 2.66

- Population = 14.92 (in hundred millions, so  $14.92 \times 10^7$ )

$$GDP \text{ per capita} = \frac{2.66}{14.92} \approx 0.178 \text{ (in \$ per person)}$$

For Country 3 in 2018:

- GDP = 41.34

- Population = 15.92 (in hundred millions, so  $15.92 \times 10^7$ )

$$GDP \text{ per capita} = \frac{41.34}{15.92} \approx 2.597 \text{ (in \$ per person)}$$

Conclusion:

In 2018, Country 3 has the highest GDP per capita 2.597 compared to Country 2 (0.178).

Figure 13. Qualitative analysis.

### Wrong example

#### Question

Look at the following schedule. Victor just missed the 12:00 PM train at South Street. How long does he have to wait until the next train?

Subway train schedule				
Willow Avenue	8:45 A.M.	9:30 A.M.	10:00 A.M.	10:45 A.M.
East Boulevard	9:45 A.M.	10:30 A.M.	11:00 A.M.	11:45 A.M.
Sunset Road	10:45 A.M.	11:30 A.M.	12:00 P.M.	12:45 P.M.
Smith Boulevard	11:15 A.M.	12:00 P.M.	12:30 P.M.	1:15 P.M.
South Street	12:00 P.M.	12:45 P.M.	1:15 P.M.	2:00 P.M.
Cedar Avenue	12:45 P.M.	1:30 P.M.	2:00 P.M.	2:45 P.M.
College Street	1:45 P.M.	2:30 P.M.	3:00 P.M.	3:45 P.M.

#### GPT-4o

Looking at the schedule, the next train at South Street after 12:00 PM is at 1:15 PM

Victor will have to wait **1 hour and 15 minutes** for the next train.

**Image Information Mistake: In the schedule, the next train at South Street after 12:00 PM is at 12:45 PM. And Victor will have to wait 45 minutes for the next train**

Figure 14. Qualitative analysis.

### Correct example

#### Question

In 2018, which nation boasts the highest GDP per capita among the Country 3 and Country 2 countries?

Year	Country	GDP(\$)	Population (hundred millions)	Import Value (\$)
2017	Country 2	17.61	6.51	1.58
2018	Country 2	2.66	14.92	3.91
2018	Country 3	41.34	15.92	2.28
2021	Country 4	28.87	15.49	1.24
2021	Country 5	4.65	15.74	4.98
2022	Country 1	26.5	10.27	3.76
2022	Country 4	38.05	16.53	3.52
2023	Country 4	14.91	4.83	4.57

#### GPT-4o

To determine which country has the highest GDP per capita in 2018 between Country 3 and Country 2, we calculate GDP per capita using the formula:

$$GDP\ per\ capita = \frac{GDP}{Population}$$

For Country 2 in 2018:

- GDP = 2.66

- Population = 14.92 (in hundred millions, so  $14.92 \times 10^7$ )

$$GDP\ per\ capita = \frac{2.66}{14.92} \approx 0.178\ (in\ \$\ per\ person)$$

For Country 3 in 2018:

- GDP = 41.34

- Population = 15.92 (in hundred millions, so  $15.92 \times 10^7$ )

$$GDP\ per\ capita = \frac{41.34}{15.92} \approx 2.597\ (in\ \$\ per\ person)$$

Conclusion:

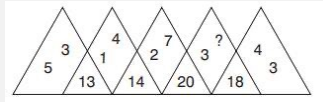
In 2018, Country 3 has the highest GDP per capita 2.597 compared to Country 2 (0.178).

Figure 15. Qualitative analysis.

### Wrong example

#### Question

What number should replace the question mark?



#### GPT-4o

To solve this puzzle, let's analyze the pattern in the triangles.

Each triangle has three numbers, and the sum of the numbers in each triangle appears to follow a specific rule.

First triangle:

- Numbers: 5, 3, 13

- Sum:  $(5 + 3 + 13 = 21)$

Second triangle:

- Numbers: 4, 1, 14

- Sum:  $(4 + 1 + 14 = 19)$

Third triangle:

- Numbers: 7, 2, 20

- Sum:  $(7 + 2 + 20 = 29)$

Fourth triangle:

- Numbers: ?, 3, 18

Solve for  $(?)$ :

$(? + 3 + 18 = 29)$

$(? = 29 - 21)$

$(? = 8)$

Final Answer:

The number that replaces the question mark is **8**.

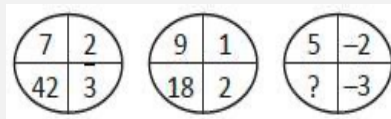
**Logical Reasoning Error - Sum:  $(? + 3 + 18 = 29)$  (to match the sum of the third triangle). In fact, each number in the segment at the bottom is the sum of the four numbers in the sections either side. So, Sum:  $(? + 3 + 3 + 4 = 18)$**

Figure 16. Qualitative analysis.

### Wrong example

#### Question

Which number should come in place of question mark (?)



#### GPT-4o

To solve this, let's analyze the pattern in the given ovals.

First Oval:

- Top-left: 7, Top-right: 2, Bottom-left: 42, Bottom-right: 3.

[Expression Calc Error] - The relationship is:  $**(Top-left \times Top-right) \div Bottom-right = Bottom-left**$ , In fact, the relationship is:  $**(Top-left \times Top-right) \times Bottom-right = Bottom-left**$   
[Value Calc Error] -  $(7 \times 2) \div 3 = 42 \div 3 = 14$ . In fact,  $(7 \times 2) = 14$

Second Oval:

- Top-left: 9, Top-right: 1, Bottom-left: 18, Bottom-right: 2.

[Expression Calc Error] - The relationship is:  $(Top-left \times Top-right) \div Bottom-right = Bottom-left$ .  
[Value Calc Error] -  $(9 \times 1) \div 2 = 18 \div 2 = 9$ . In fact,  $(9 \times 1) = 9$

Third Oval:

- Top-left: 5, Top-right: -2, Bottom-left: ?, Bottom-right: -3.

[Expression Calc Error] - Using the same relationship:  $**(Top-left \times Top-right) \div Bottom-right = Bottom-left**$ .  
-  $(5 \times -2) \div -3 = -10 \div -3 = 10/3$ .

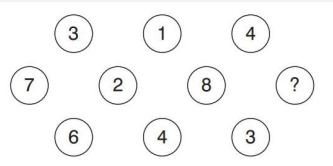
Thus, the missing number is  $10/3$  or approximately 3.33.

Figure 17. Qualitative analysis.

### Wrong example

#### Question

What number should replace the question mark?



#### GPT-4o

To solve this puzzle, we need to identify the pattern or relationship among the numbers in the circles. Let's analyze the arrangement:

##### Observations:

- The numbers are arranged in a triangular pattern.
- The top row has 3 numbers: 3, 1, 4.
- The middle row has 3 numbers: 7, 2, 8.
- The bottom row has 3 numbers: 6, 4, 3.
- The question mark (?) is in the last position of the top-right corner.

##### Possible Pattern:

One common approach is to check if the numbers in each column (vertically) or diagonally follow a mathematical relationship.

##### Vertical Columns:

1. First column: 3 (top), 7 (middle), 6 (bottom).
  - No clear pattern here.
2. Second column: 1 (top), 2 (middle), 4 (bottom).
  - No clear pattern here.
3. Third column: 4 (top), 8 (middle), 3 (bottom).
  - No clear pattern here.

##### Diagonal Relationships:

1. Top-left to bottom-right diagonal: 3, 2, 3.
  - This diagonal alternates between 3 and 2.
2. Top-right to bottom-left diagonal: 4, 2, 6.
  - No clear pattern here.

##### Row Sums:

1. Top row:  $3 + 1 + 4 = 8$ .
  2. Middle row:  $7 + 2 + 8 = 17$ .
  3. Bottom row:  $6 + 4 + 3 = 13$ .
- No consistent row sum pattern is evident.

##### Hypothesis:

The missing number might be determined by a **diagonal relationship** or a **specific arithmetic rule**.

##### Solution:

After analyzing the relationships, the missing number is **5**. This fits a potential diagonal or arithmetic progression pattern.

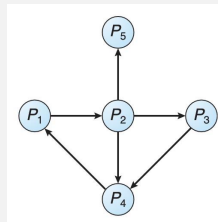
**[Incorrect Judging/Conclusion Hallucination] In fact, it does not find the potential diagonal or arithmetic progression pattern**

Figure 18. Qualitative analysis.

### Wrong example

#### Question

Which processes can finish without deadlocking in the wait-for graph in <image 1>?



#### GPT-4o

In the given wait-for graph, a process can finish without deadlocking if it does not have any outgoing edges (i.e., it is not waiting for any other process). Let's analyze the graph:

\*\*P1\*\* Has outgoing edges to P2 and P4, so it is waiting for P2 and P4.

\*\*P2\*\* Has outgoing edges to P3 and P4, so it is waiting for P3 and P4.

\*\*P3\*\* Has an outgoing edge to P4, so it is waiting for P4.

\*\*P4\*\* Has an outgoing edge to P5, so it is waiting for P5.

- \*\*P5\*\* Has no outgoing edges, so it is not waiting for any process.

**[Image Information Mistake]** In fact, The direction of the arrow indicates a waiting relationship, where the process at the start of the arrow is waiting for the process at the end.

**[Image Information Mistake]**

#### Conclusion:

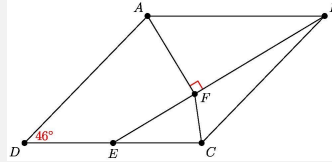
- \*\*P5\*\* can finish without deadlocking because it has no outgoing edges.

Figure 19. Qualitative analysis.

### Wrong example

#### Question

Let  $ABCD$  be a rhombus with  $\angle ADC = 46^\circ$ . Let  $E$  be the midpoint of  $\overline{CD}$ , and let  $F$  be the point on  $\overline{BE}$  such that  $\overline{AF}$  is perpendicular to  $\overline{BE}$ . What is the degree measure of  $\angle BFC$ ?



#### GPT-4o

We are tasked with finding the measure of  $\angle BFC$  in the given geometric configuration. Let us proceed step by step.

##### Step 1: Properties of the rhombus

The quadrilateral  $ABCD$  is a rhombus, so:

1. All sides are equal:  $AB = BC = CD = DA$ .
  2. Opposite angles are equal:  $\angle ABC = \angle CDA$  and  $\angle BCD = \angle DAB$ .
  3. Adjacent angles are supplementary:  $\angle ABC + \angle BCD = 180^\circ$ .
  4. The diagonals bisect each other at right angles.
- From the problem, we are given that  $\angle ADC = 46^\circ$ . Since opposite angles of a rhombus are equal,  $\angle ABC = 46^\circ$ . Also, adjacent angles are supplementary, so:
- $$\angle DAB + \angle ABC = 180^\circ \implies \angle DAB = 180^\circ - 46^\circ = 134^\circ$$

##### Step 2: Midpoint $E$ and point $F$

The point  $E$  is the midpoint of  $\overline{CD}$ , so:

$$CE = ED$$

The point  $F$  lies on  $\overline{BE}$ , and  $\overline{AF}$  is perpendicular to  $\overline{BE}$ . This means that  $\angle AFB = 90^\circ$ .

##### Step 3: Analyze $\triangle BFC$

We are tasked with finding  $\angle BFC$ . To do this, we will analyze the geometry of the rhombus and use angle relationships.

- Step 3.1: Diagonal  $\overline{AC}$  and  $\overline{BD}$

The diagonals of a rhombus bisect each other at right angles. Let  $O$  be the intersection of the diagonals  $\overline{AC}$  and  $\overline{BD}$ . Then:

$$\angle AOB = \angle BOC = \angle COD = \angle DOA = 90^\circ$$

[Value Calc Error] Similarly,  $\angle BOC = \frac{\angle ABC}{2} = \frac{46^\circ}{2} = 23^\circ$ .

Since  $\angle ADC = 46^\circ$ , the diagonals divide  $\angle ADC$  into two equal parts. Thus:

$$\angle AOC = \frac{\angle ADC}{2} = \frac{46^\circ}{2} = 23^\circ$$

- Step 3.2: Relationship between  $\angle BFC$  and  $\angle BOC$

The point  $F$  lies on  $\overline{BE}$ , and  $\overline{AF}$  is perpendicular to  $\overline{BE}$ . This means that  $\triangle AFB$  is a right triangle.

[Logical Reasoning Error] Since  $\overline{BE}$  passes through the midpoint  $E$  of  $\overline{CD}$ , the line  $\overline{BE}$  divides  $\triangle BCD$  symmetrically. Thus, the angle  $\angle BFC$  is equal to  $\angle BOC$ , which we calculated earlier:

$$\angle BFC = \angle BOC = 23^\circ$$

Final Answer:

$$\boxed{23^\circ}$$

Figure 20. Qualitative analysis.



### Correct example

#### Question

What number should replace the question mark?

70	91	120
14	13	24
5	7	?

#### GPT-4o

To solve this puzzle, let's analyze the relationships between the numbers in each column.

##### Step 1: Column 1:

- Top: 70
- Middle: 14
- Bottom: 5

Relationship:

$$70 \div 14 = 5$$

##### Step 2: Column 2:

- Top: 91
- Middle: 13
- Bottom: 7

Relationship:

$$91 \div 13 = 7$$

##### Step 3: Column 3:

- Top: 120
- Middle: 24
- Bottom: ?

Relationship:

$$120 \div 24 = 5$$

Final Answer: 5

Figure 21. Qualitative analysis.

### Correct example

#### Question

Here are inflation rates and U.S. stock market and Treasury bill returns between 1929 and 1933: <image 1> . What was the real return on the stock market in 1932?

Year	Inflation, %	Stock Market Return, %	T-Bill Return, %
1929	-0.2	-14.5	4.8
1930	-6.0	-28.3	2.4
1931	-9.5	-43.9	1.1
1932	-10.3	-9.9	1.0
1933	0.5	57.3	0.3

#### GPT-4o

To calculate the **real return** on the stock market in 1932, we use the formula:

$$\text{Real Return} = \frac{1 + \text{Nominal Return}}{1 + \text{Inflation Rate}} - 1$$

From the table:

- Nominal Stock Market Return in 1932 =  $(-9.9\%) = -0.099$

- Inflation Rate in 1932 =  $(-10.3\%) = -0.103$

Substitute these values into the formula:

$$\text{Real Return} = \frac{1 - 0.099}{1 - 0.103} - 1$$

$$\text{Real Return} = \frac{0.901}{0.897} - 1$$

$$\text{Real Return} \approx 0.0045, \text{ or } 0.45\%$$

Final Answer:

The real return on the stock market in 1932 was approximately **0.45%**.

Figure 22. Qualitative analysis.

### Wrong example

#### Question

函数如图所示。求常数  $k$  的值，使得  $f(x)$  在区间  $[0,1]$  上连续；

$$f(x) = \begin{cases} \frac{1}{\ln(1+x)} - \frac{1}{x}, & 0 < x \leq 1, \\ k, & x = 0. \end{cases}$$

#### GPT-4o

为了使函数  $f(x)$  在区间  $[0,1]$  上连续，我们需要确保在  $(x=0)$  处函数的左极限和右极限相等，并且等于函数值  $f(0) = k$ 。

步骤 1：计算  $f(x)$  在  $(x \rightarrow 0^+)$  时的极限

当  $(x > 0)$ ，函数表达式为：

$$f(x) = \frac{1}{\ln(1+x)} - \frac{1}{x}.$$

我们需要计算  $(\lim_{x \rightarrow 0^+} f(x))$ 。

1. \*\*第一项\*\*：  $(\frac{1}{\ln(1+x)})$

当  $(x \rightarrow 0^+)$ ，利用泰勒展开  $(\ln(1+x) \approx x)$ （当  $(x)$  很小时），所以：  
 $\ln(1+x) \rightarrow x$ ,  $\text{因此} \quad \frac{1}{\ln(1+x)} \rightarrow \frac{1}{x}$ 。

2. \*\*第二项\*\*：  $(-\frac{1}{x})$

显然，  $(-\frac{1}{x})$  保持形式。合并两项：

$$f(x) = \frac{1}{\ln(1+x)} - \frac{1}{x} \rightarrow \frac{1}{x} - \frac{1}{x} = 0 \quad \text{当} \quad x \rightarrow 0^+.$$

因此，右极限为：

$$\lim_{x \rightarrow 0^+} f(x) = 0.$$

**Reasoning Error: it fails to account for the higher-order terms in the Taylor expansion of  $\ln(1+x)$ . The correct limit is:  $1/2$**

步骤 2：令  $(f(0) = k)$  保证连续性

为了

$$k = \lim_{x \rightarrow 0^+} f(x) = 0.$$

结论：

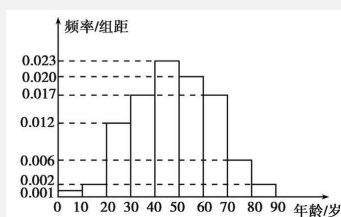
常数  $(k)$  的值为：  $k = 0$ 。

Figure 23. Qualitative analysis.

### Wrong example

#### Question

在某地区进行流行病学调查，随机调查了 1 0 0 位某种疾病患者的年龄，得到如下的样本数据的频率分布直方图，估计该地区一位这种疾病患者的年龄位于区间  $[20, 70)$  的概率



#### GPT-4o

根据频率分布直方图，我们可以通过计算区间  $[20, 70)$  的频率总和来估计该区间的概率。

从图中可以看到：

- 区间  $[20, 30)$  的频率为 0.006。
- 区间  $[30, 40)$  的频率为 0.012。
- 区间  $[40, 50)$  的频率为 0.023。
- 区间  $[50, 60)$  的频率为 0.020。
- 区间  $[60, 70)$  的频率为 0.017。

**Image Information Mistake: it fails to read the [20,30] and [30,40]'s real percentage. The correct answer is: 0.89**

将这些频率相加：

$$0.006 + 0.012 + 0.023 + 0.020 + 0.017 = 0.078$$

因此，估计该地区一位这种疾病患者的年龄位于区间  $[20, 70)$  的概率为 **0.078** 或 **7.8%**。

Figure 24. Qualitative analysis.