



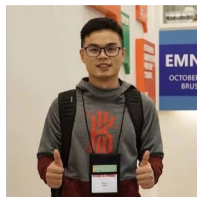
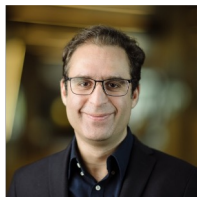
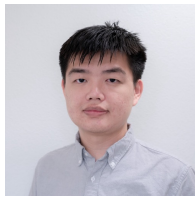
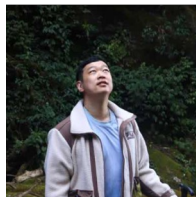
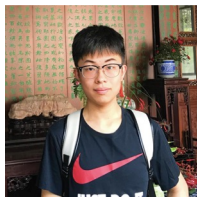
FIRE: A Dataset for Feedback Integration and Refinement Evaluation of Multimodal Models

Pengxiang Li^{*1,2}, Zhi Gao^{*1,3}, Bofei Zhang^{*1}, Tao Yuan¹, Yuwei Wu^{E2},
Mehrtash Harandi⁴, Yunde Jia¹, Song-Chun Zhu^{1,3,5}, Qing Li^{E1}

¹ BIGAI ² BIT ³ Peking University ⁴ Monash University ⁵ Tsinghua University

* Equal contribution, ^E Corresponding Author

<https://mm-fire.github.io>



Motivation

We tackle the limitation of open-source Vision-Language Models (VLMs) in feedback-refining within a single session, aiming to enhance their interactive learning capabilities.



Ground truth: Car.



User

Which kind of vehicle is waiting for the traffic light?

The vehicles look like typical vehicles used for personal transport. The vehicles are **motorcycles**.

Wrong



LLaVA-Next-8B



User

Motorcycles are not correct. Look again at the structure and features of the vehicles, such as **the number of wheels and the body shape.**

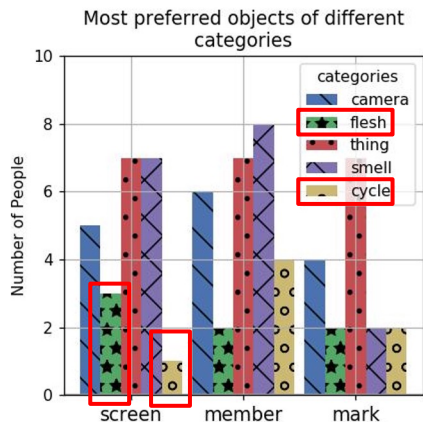
Upon reevaluating the image, the vehicles are indeed **motorcycles**.



LLaVA-Next-8B

Failed in Refining

Limited feedback-refining ability



Ground truth: No.



User

Is the object mark in the category cycle preferred by more people than the object screen in the category flesh?

Yes the object mark in the category cycle is preferred by more people than the object screen in the category flesh. The bar representing the object mark in the **cycle category** reaches up to the number **8** on the y-axis, while the bar for the object screen in the **flesh category** is just above the number **2** on the y-axis.



Mini-Gemini-34B-HD



User

Your answer is **incorrect**. You correctly identified the need to compare the number of people, but your conclusion is wrong. Look closely at the bar heights for the **object mark in the cycle category** and the **object screen in the flesh category**. Try again.


I apologize for the mistake in my previous answer. Upon closer inspection, the object mark in the cycle category is indeed preferred by more people than the object screen in the flesh category. The bar for the object mark in the **cycle category** reaches up to the number **8** on the y-axis, while the bar for the object screen in the **flesh category** is just above the number **2** on the y-axis. Therefore, the correct answer is **yes**, the object mark in the cycle category is preferred by more people than the object screen in the flesh category.



Mini-Gemini-34B-HD

**Same as initial response.
Failed in Refining**

Limited feedback-refining ability



Ground truth:
Nutritional deficiency.

User: Why are the leaves turn yellow?

QwenVL-72B: Overwatering

User: Your answer is incorrect. Please judge the possible reasons based on the texture characteristics of the yellowing leaves. Please note that overwatering will not cause the leaves to turn yellow in such pattern. Yellowing leaves may be caused by a disease in the plant.

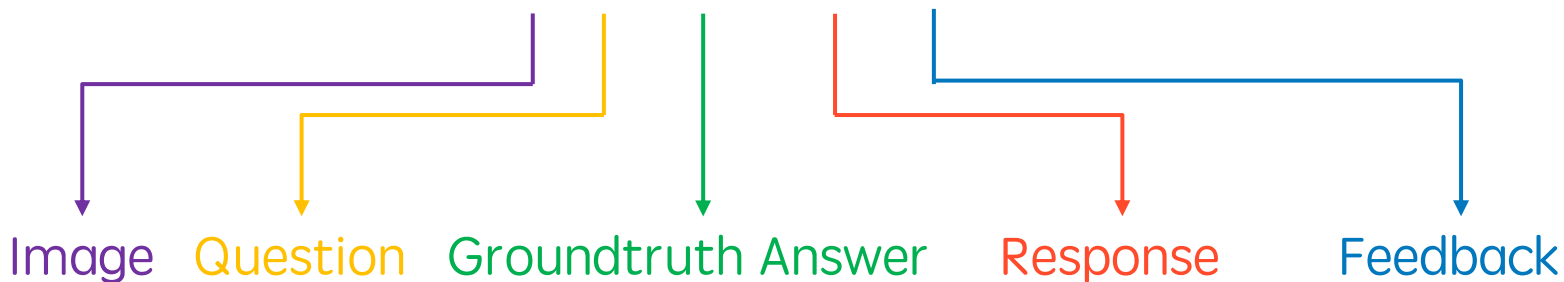
QwenVL-72B: Overwatering

The image shows a chat interface with a user asking a question about a plant with yellowing leaves. The ground truth is nutritional deficiency. The model's first response is "Overwatering", which is marked as "Wrong". The user provides feedback, stating that overwatering is incorrect and that the yellowing pattern is due to a disease. The model's second response is also "Overwatering", which is marked as "Failed in Refining".

One reason of such limitation is that instruction-following datasets predominantly consist of data focused on direct task completion, with **little to no inclusion** of feedback-refining information.

Task setting

Data format

$$\{I, q, gt, \{r_i, f_i\}_1^n\}$$


Dialogue simulation



Vision Language Model
ROLE: Teacher
Playing: Scoring and feedback.

Input: $I, q, \textcircled{gt}, r_{i-1}$
Output: f_i

Vision Language Model
ROLE: Student
Playing: Answer the Question.

Input: I, q, f_{i-1}
Output: r_{i+1}



Groundtruth Answer is only visible to teacher model.

Dataset statistic



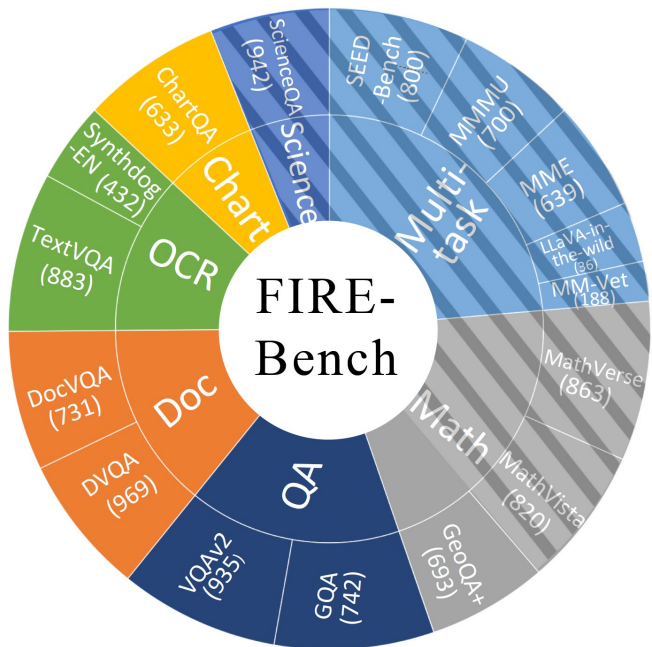
Training Set FIRE-100K & FIRE-1M



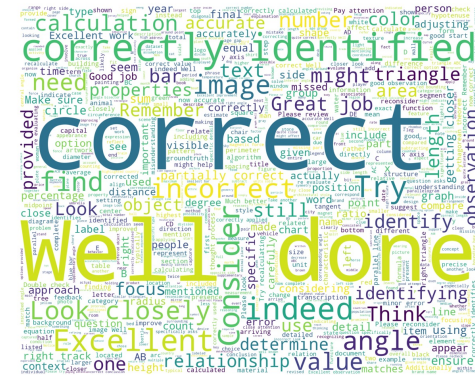
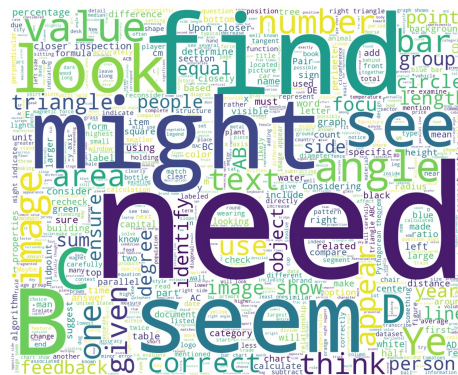
Statistic	FIRE-100K	FIRE-1M
Total dialogue	104941	1003490
Avg. Response Length	41.20	79.33
Avg. Feedback Length	41.74	54.87

FIRE offers **1.1 million** high-quality feedback-refining dialogues across **8** diverse tasks and **18** various data sources.

Dataset statistic



Test Set FIRE-Bench

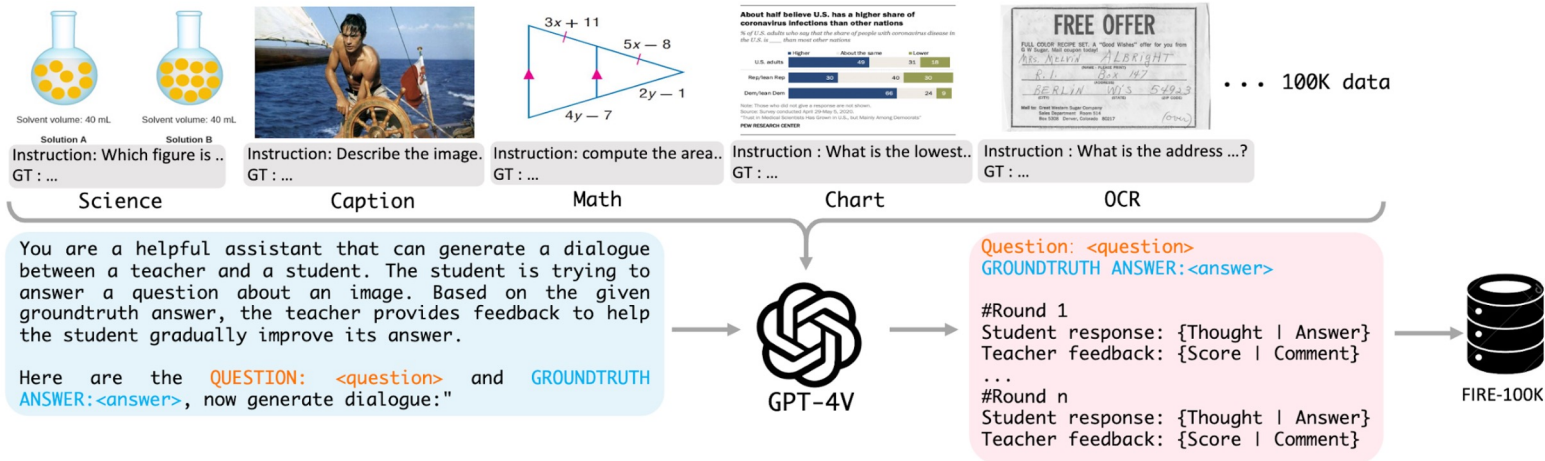


Statistic	FIRE-Bench
Total dialogue	11006
Avg. Response Length	38.70
Avg. Feedback Length	41.74

FIRE-Bench offers **11 K** high-quality feedback-refining dialogues across **8** diverse tasks and **16** various data sources, including 8 test splits from the training data sources and **8 completely new benchmarks**.

Data collection

(a) FIRE-100K Data Generation



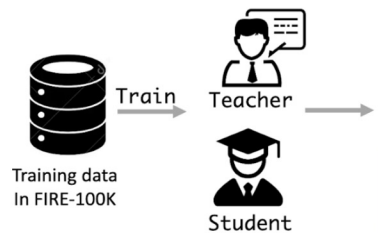
We use GPT-4V to generate 100K data (FIRE-100K), with each dialogue was generated in a one-go manner.

Data collection

We use FIRE-100K to train a Teacher model and a Student model.

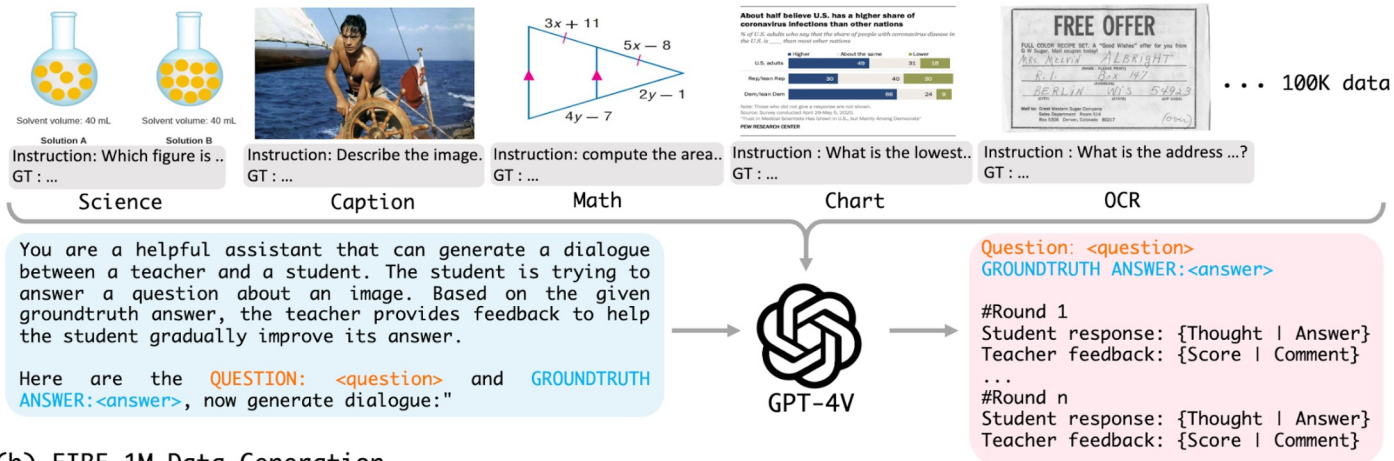
Then we simulate dialogues between the student and the teacher models and scale up the data points to ~1M, rendering a split named FIRE-1M.

(b) FIRE-1M Data Generation

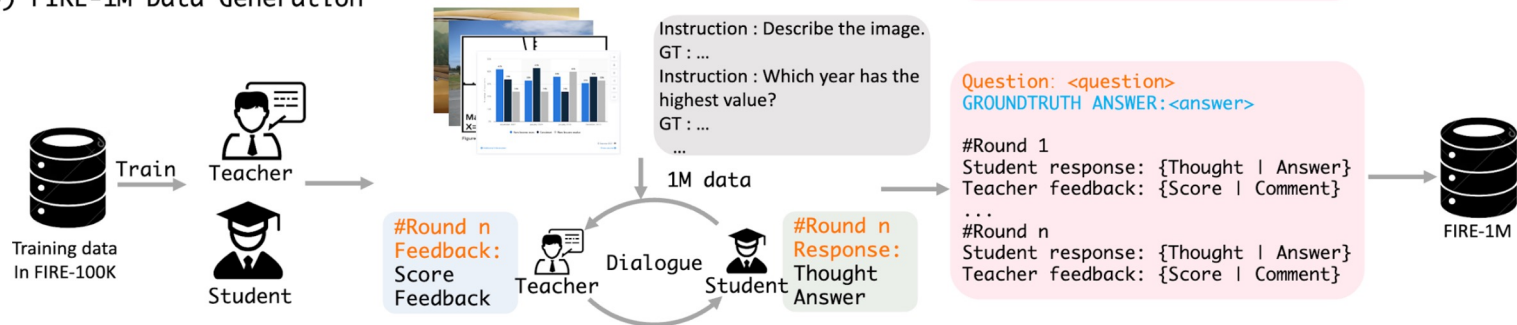


Data collection

(a) FIRE-100K Data Generation



(b) FIRE-1M Data Generation



Training-Student Model

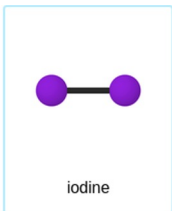
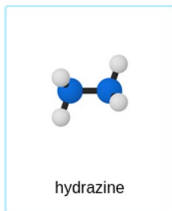
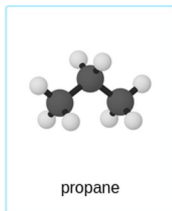
$$\{I, q, gt, \{r_i, f_i\}_1^n\}$$

Question:

Look at the models of molecules below. Select the elementary substance.

Choices:

A. propane B. hydrazine C. iodine



The elementary substance in the image is represented by the molecule labeled "propane." Propane is a simple organic molecule composed of three carbon atoms and one hydrogen atom. It is a gas at room temperature and is commonly used as a fuel. The other molecules, hydrazine and iodine, are more complex compounds. Hydrazine is a compound of hydrogen and nitrogen, while iodine is a non-metallic element. Therefore, the correct answer to the question is: **A. propane**



Score: 3.

Comments: You correctly identified that propane is a simple organic molecule composed of carbon and hydrogen atoms. However, the question asks for the elementary substance, which refers to a basic chemical element. Propane is a compound, not an element. Consider the other options and think about which one is an element.



Text completion:



Ground-truth Answer: 🐶

$$\min \mathbb{E}_{(I, q, gt, \{r_i, f_i\}_{i=1}^n) \sim \mathbb{D}} \left[- \sum_{i=2}^n \log P(r_i | I, q, \{r_j, f_j\}_{j=1}^{i-1}) \right]$$

Training-Teacher Model

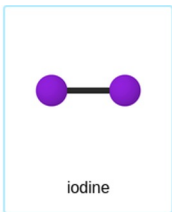
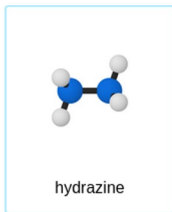
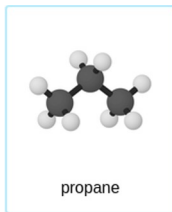
 $\{I, q, gt, \{r_i, f_i\}_1^n\}$

Question:

Look at the models of molecules below. Select the elementary substance.

Choices:

A. propane B. hydrazine C. iodine



I apologize for the confusion. You are correct that propane is a compound, not an element. The elementary substances in the image are represented by the molecules labeled "hydrazine" and "iodine". **Hydrazine is a compound of hydrogen and nitrogen, and iodine is a non-metallic element.** Therefore, the correct answer to the question is: **C. iodine**



Text completion:

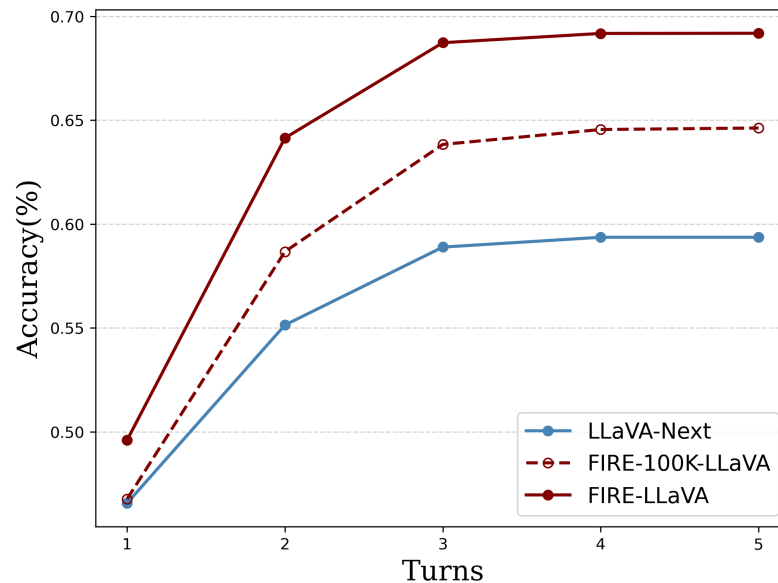
Ground-truth Answer:

C. iodine

$$\min \mathbb{E}_{(I, q, gt, \{r_i, f_i\}_{i=1}^n) \sim \mathbb{D}} \left[- \sum_{i=1}^n \log P(f_i | I, q, gt, \{r_j, f_j\}_{j=1}^{i-1}, r_i) \right]$$

Experiments

LLaVA-Next fine-tuned by FIRE can complete the response **faster** and has **better feedback refining capabilities** (7% more improvements) under the same feedback model.



Experiments

Average Turn (AT). The AT metric evaluates **how fast** a VLM could achieve a satisfactory result based on feedback.

Average Dialogue Refinement (ADR). The ADR metric evaluates how much knowledge VLMs could learn from feedback **in a dialogue**.

Average Turn Refinement (ATR). ATR evaluates how much knowledge VLMs could learn from feedback **in one turn**.

Refinement Ratio (RR). RR measures the proportion of data that have a wrong initial response and a correct final response

Results in free dialogue over all test data in FIRE.

Model	AT (↓)	ADR (↑)	ATR (↑)	RR (↑)
LLaVA-Next-8B	1	0.97	0.41	0.25
FIRE100K-LLaVA-8B	0.92	1.27	0.55	0.34
FIRE-LLaVA-8B	0.84	1.56	0.66	0.39

LLaVA-Next fine-tuned by FIRE can complete the response **faster** and has **better feedback refining capabilities**.

Visualization

Question:

What is the color of the gloves worn by one of the football players in the image?

- A. White
- B. Black
- C. Green
- D. Blue



Ground-truth Answer:

B. Black



The color of the gloves worn by the football player in the image is white. This can be determined by observing the player in the white uniform, who is on top of the other player. The gloves match the color of his uniform, which is predominantly **white** with blue detailing. Therefore, the correct answer to the question is: **A. White**



Score: 4.

Comments: You correctly identified that the gloves match the color of the player's uniform, which is predominantly white with blue detailing. However, the question specifically asks about **the color of the gloves** worn by one of the football players, **not the uniform**. **Look closely at the gloves worn by the player in the green uniform, who is on the ground.** Consider the options provided and re-evaluate your answer.



Upon re-evaluation, the gloves worn by the player in the green uniform, who is on the ground, are black. This can be confirmed by the contrast between the gloves and the green uniform, as well as the white uniform of the player on top. Therefore, the correct answer to the question is: **B. Black**



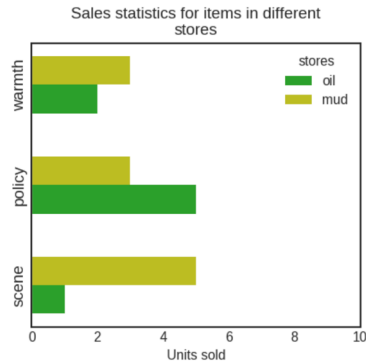
Score: 10.

Comments: Excellent observation! You correctly identified that the gloves worn by the player in the green uniform are black. This matches the ground truth answer. Well done on adjusting your focus and re-evaluating the image.

Visualization

Question:

How many units of the item scene were sold in the store oil?



Based on the image, it is not possible to determine the exact number of units of the item "scene" sold in the store "oil" because the bar representing "oil" for "scene" is not fully visible. Only a portion of the bar is shown, and without the full length of the bar or a scale, the number of units sold cannot be accurately determined from the given image.



Score: 3.

Comments: You are correct that the image does not provide enough information to determine the exact number of units of the item "scene" sold in the store "oil." However, you should **consider the visible portion of the bar and any patterns or trends that might help estimate the number of units sold.** Try to use the visible information to make an educated guess.

Upon closer inspection, the visible portion of the "oil" bar for "scene" appears to extend slightly beyond the halfway point between 0 and 2 on the x-axis. This suggests that a small number of units were sold, possibly between 1 and 2. However, without more information or a clearer view of the entire bar, it is not possible to provide an exact number.



Ground-truth Answer:

1



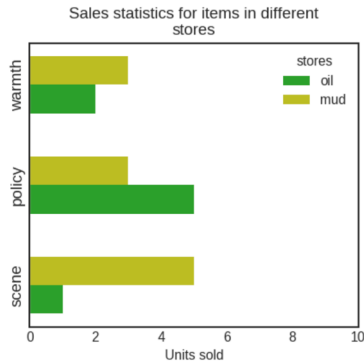
Score: 6.

Comments: Your observation about the visible portion of the "oil" bar for "scene" is a good start. You correctly noted that the bar extends slightly beyond the halfway point between 0 and 2. This suggests that a small number of units were sold. However, to improve your answer, **consider the possibility of the bar extending just beyond the halfway point and estimate the number of units sold more precisely.** Remember, the exact number cannot be determined without more information, but you **can make an educated guess based on the visible data.**

Visualization

Question:

How many units of the item scene were sold in the store oil?



Ground-truth Answer:

1

Considering the visible portion of the "oil" bar for "scene" extends just beyond the halfway point between 0 and 2, **it is reasonable to estimate that approximately 1 unit of the item "scene" was sold in the store "oil."** This is based on the assumption that the bar extends equally beyond the halfway point on both sides, which is a common pattern for bars in such graphs. However, without more information or a clearer view of the entire bar, this remains an educated guess rather than a definitive answer.



Score: 9.

Comments: Considering the visible portion of the "oil" bar for "scene" extends just beyond the halfway point between 0 and 2, **it is reasonable to estimate that approximately 1 unit of the item "scene" was sold in the store "oil."** This is based on the assumption that the bar extends equally beyond the halfway point on both sides, which is a common pattern for bars in such graphs. However, without more information or a clearer view of the entire bar, this remains an educated guess rather than a definitive answer.



Takeaway

We present **FIRE**, a feedback-refinement dataset with **1.1M multi-turn conversations**, which empowers VLMs to refine their responses based on given feedback.

We develop **FIRE-LLaVA**, an advanced VLM that could improve its responses based on feedback, making efficient interaction between users and VLMs.

We have developed and released a **Feedback Model** designed to enhance VLM performance by providing targeted feedback.

Code & Dataset: <https://mm-fire.github.io>



FIRE: A Dataset for Feedback Integration and Refinement Evaluation of Multimodal Models

